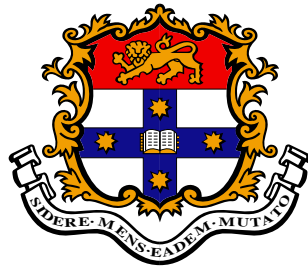


Family: Documents, files, nodes and relationships

Sam Thorogood

CHAI Group, School of IT
The University of Sydney
NSW 2006 Australia

sam.thorogood@student.usyd.edu.au



November 7, 2008

Abstract

Existing file systems fail to support users in identifying important relationships between their own documents. For example, when a user copies a file on their computer, their mental model may link the copy with the original. However, modern operating systems treat these as unrelated. Moreover, it is not straightforward to create a system that can recognise these kinds of relationships.

This thesis describes Family, a system which aims to support users in seeing the temporal sequence of files that are part of the history, or lineage, of any nominated file. We term these files the *ancestors* of this nominated file, and the inverse relationship the *descendants*. We also identify key conceptual models that constitute the cases of important family relationships. These correspond to simple operations on files (such as copy, rename) that may be implemented in various different ways at the operating system level, depending upon the application used to undertake these operations. These are operations that will typically involve multiple complex actions at an operating system level, but which are conceptually simple to end-users.

We introduce a formal algorithm for Family, as well as its specific implementation, known as Lighthouse. This algorithm is defined explicitly, describing how the broad concepts of *ancestors* and *descendants* are found through low-level analysis and observation of low-level system activity. Finally, we evaluate Family and Lighthouse in terms of correctness versus a predefined standard, as well as performing a short performance evaluation.

Acknowledgements

Firstly, I would like to thank the creators of UNIX for popularising the wonderful modern file system that all computer users must interact with today. Secondly, I must thank Apple for creating a wonderful environment on which to develop and obtain the most relevant of information for use within my implementation.

This thesis would not have been possible without the wonderful support of my academic group here at the University of Sydney. I would like to thank Judy Kay, my supervisor, for her support and suggestions. My fellow Honours students Greg and James were instrumental in sanity-checking all my code. Coffee runs with Greg, Mat and Katie have also kept me sane. Tim Tams on Thursdays (or was it Tuesdays) kept me going with regular chocolate injections – thanks to Trent for originally organising this. Not to mention your wisdom and support throughout the first half of the year. I would also like to thank Anthony, a PhD student, for my incessant questions about LaTeX. Also worthy of mention are my fellow Honours students – again, Greg and James, along with Mat, Katie, and others – I hope to keep in touch with all of you for many years to come.

Of course, I have to mention my family – Mary, Peter, Jen and Jay. They have put up with my late nights and late mornings, making sure I still eat and keep healthy, as well as providing feedback and editing services.

Finally, many thanks to the Apple University Consortium in Australia for their support through both an Honours scholarship and an Innovation Seeding Grant.

Contents

1	Introduction	1
1.1	Thesis Structure	4
1.2	Vision & Motivation	5
1.3	Challenges	6
1.4	Contributions	7
2	Related Work	9
2.1	Files	10
2.2	Layout	11
2.3	Temporal Usage & Relationships	11
2.3.1	Explicit	12
2.3.2	Implicit	12
3	Approach	13
3.1	High-level Formalism	15
3.2	Conceptual Models	15
3.3	Low-level analysis	17
3.4	Algorithm	19
4	Implementation	21
4.1	Kernel Extension	21
4.2	Lighthouse	22
4.3	Testing Framework	22
4.4	Walkthrough	24
5	Evaluation	25
5.1	Correctness	25
5.2	Performance	27
6	Conclusion	31
6.1	Contributions	31
6.2	Future Work	31
6.3	Conclusions	32
A	Kernel Extension	35
B	Lighthouse	52

List of Figures

1.1	A simple relationship example	2
1.2	Temporal relationship example	4
2.1	A file system tree	10
3.1	Concept leap	14
3.2	Identifying a copied file	19
4.1	Mapping files to nodes	23
4.2	Walkthrough of the testing framework	24
5.1	Single case test results	28
5.2	Timing for traditional file system actions	29
5.3	Time taken for correct mapping	30

Chapter 1

Introduction

Every day, users undertake a multitude of conceptual actions over their own files in the pursuit of achieving their computing goals. These actions are typically well-known to users and are well-defined by most applications that work with documents. For example: a user might *copy* a report template, before editing it and performing a *save*, before *moving* the file to a USB stick.

As part of this daily process, a user may build up a *mental model* of how their files relate to each other, based on these simple conceptual actions. As per our example, the end-user may hold a simple understanding that the file they have just worked on is based on the original template. However, modern operating systems are not able to capture this seemingly simple relationship, and others like it. From the point of view of Mac OS X, Windows, and Linux, these two files are independent and unrelated sequences of raw data. One of the key reasons for this is that the seemingly simple conceptual actions previously mentioned, including *copy* and *save*, are typically implemented as a multitude of lower-level system actions.

We introduce an example in Figure 1.1, and describe it below.

John is working on a conference report. Initially, he stores his notes in a file named **Notes**. Later, he creates a new file, **Final Report**, based on his notes. At this point, he may reasonably associate the two documents within a mental model of his personal dataspace.

As noted above, modern operating systems will treat these two new documents as unrelated – as from a functional point of view, these are now two completely unrelated files. However, such relationships can be important for finding the right file for a particular task, or when an activity calls for key related files that were used to create an existing file. John may therefore find it useful to have this information available to him, rather than relying on his own personal mental model, as continued below.

Days later, John is working on **Final Report**. He realises that an important figure is missing from his report – he may have deleted it accidentally. However, he cannot remember exactly which notes he originally used to create this document, as they have already been filed away into another directory which contains presentation notes from all conferences held in 2008.

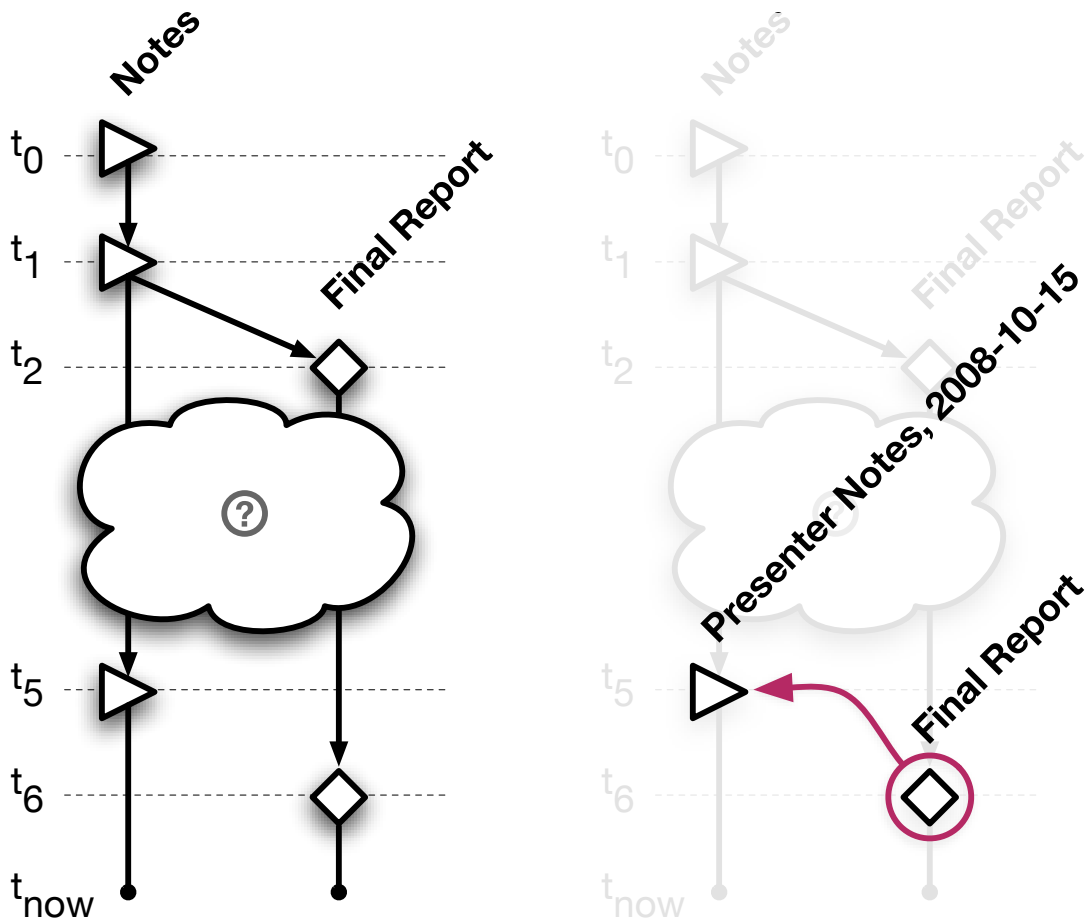


Figure 1.1: A simple relationship example, represented by the left graph. Each icon represents a unique state of a file – for example, **Notes** was created at t_0 , and updated at t_1 and t_5 . **Final Report** was copied from **Notes** at t_2 , and updated at t_6 . After some time (represented by the cloud), the two files may appear unrelated, but they retain a common ancestry. Family attempts to maintain information about this ancestry, and it could be used to simplify the relationship between the ‘current’ versions of these two files – as shown in the right graph – to be presented to an end-user.

John uses Family to find files that are related to **Final Report**, revealing his original set of notes. The document is now named **Presenter Notes, 2008-10-15** and is located in his presentation notes directory.

Relationship Classes

Broadly, we are concerned with two classes of relationships between files. The first, as described through the above example, relates files through conceptual actions such as *copy*, *duplicate*, or *cut and paste*. As mentioned previously, modern operating systems are not able to capture relationships of this type.

The second class of relationship links the unique states of a file over time. As files are changed, modified and *saved* over time, their content will exist in various different states. However, within most modern file systems – such as HFS+, NTFS, ext3 and XFS – files are mutable. While these file systems can provide end-users with the time each file was *last* changed, there is typically no support for retrieving information about *every* time they have been changed. In other words, there is traditionally no file system support to identify each unique state of a file over time.

There are two key motivators for being concerned about the states of a file over time. The first is a task-based goal, with the aim of relating work undertaken at a similar point in time. This information may be useful for task management or backtracking over the lifetime of one's own documents. Additionally, having a view of changes over all files – as opposed to just a single file – may help to better represent a temporal link between changes performed on many different files during a common timeframe.

We give an example of our first, task-based goal below. This is displayed in Figure 1.2.

A user, Johan, may be working on a poster within some graphic design software. He opens the file and notices a small mistake – fixing it and instinctively saving the document. He is now interested to determine when he last worked on this document, but the 'last changed' field now reflects the time at which he fixed the typo. He was interested in finding the notes he was working on around the same time as the poster, based on their temporal usage, but will now have to approximate this information.

So, although files may have been changed at similar times as part of one activity, if some of them are later changed there is no way to detect the temporal link – or *relationship* – suggested by the work done on them at the same time.

The second motivator for interest in the states of a file over time is tied to our first class of relationship – where files are related through conceptual actions such as *copy*, *save as* or *move*. Temporal relationships are important. These states are important, as the first class of relationship may relate directly to a specific state of a file, rather than to the file as a whole. For example, one file may be *copied*, before it is updated and *saved*. The duplicate file resulting from the *copy* operation is not related to the new, updated version: it is related only to the previous state of the file. We describe this relationship as forming the history of any given file.

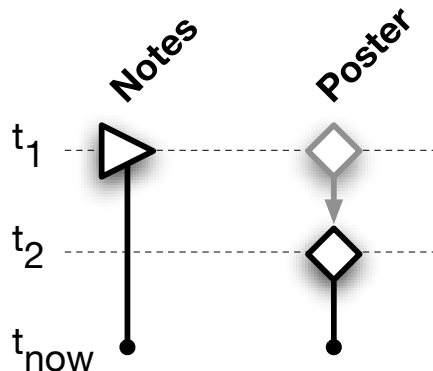


Figure 1.2: A simple example demonstrating the relevance of temporal relationships. As modern file systems only present the *last* time a file was changed, the temporal indication that these files are related (at t_1) will be lost if one of the files is changed again later (at t_2).

The combination of these two classes of relationship aims to approximate the *mental model* of how an end-user may perceive their files over time.

1.1 Thesis Structure

This thesis describes Family, an architecture which aims to support users by capturing and revealing these seemingly simple relationships, providing information about the temporal sequence of files that are part of the *lineage* of any nominated file. We also identify key conceptual actions that map directly to these important relationships which we aim to identify. We term these mappings as *conceptual models*, and a later section of this thesis aims to define them explicitly – including descriptions of how different applications may enact them in terms of lower-level system actions.

In Chapter 2, we introduce the field of “Personal Information Management”, or PIM, which has been defined as “the practice and study of the activities people perform to acquire, organize, maintain, and retrieve information for everyday use.” [1]. While the contributions of Family are very specific, we do have a more general objective of enriching this field as a whole, and this has motivated our research. Furthermore, this chapter reviews and introduces other relevant work that has a similar goal of supporting end-users in this field.

Chapter 3 details our approach to both the algorithm and architecture of Family. To find and capture relationships from conceptual actions between user files, we rely on a fairly complex algorithm. This algorithm has been implemented – something we describe in Chapter 4 – in Python, and while it requires some unique file system features (such as extended attributes) it is largely platform-independent. However, there is still a large proportion of platform specific-code which has the sole aim of providing all low-level system information, which is typically exposed in different ways on different platforms. In the case of our implementation, we target the Mac OS X platform, with this code implemented in C.

Our evaluation, both design and results, are covered in Chapter 5. We primarily focus this evaluation around determining the correctness of our algorithm.

Through our research into this area, and our implementation of Family, we have reached several core conclusions which have implications for future work. We discuss both of these elements in Chapter 6.

1.2 Vision & Motivation

Currently, PIM is a decentralised and disorganised environment. Indeed, it has previously been described as ‘fragmented’ [16]. We propose that this is fair assessment: modern computer users are regularly presented with a myriad of different applications and file formats which may all be equally as useful as each other. This leads to substantial integration difficulties.

Modern PIM is now also more decentralised than ever before. With the mainstream adoption of the internet, and the introduction of mobile computing through devices such as laptops, phones, and even the now-humble USB stick, users have been empowered to move, send, email, or even upload their data to have it accessible anywhere they are. However, this benefit came with a cost – the reality of being empowered with one’s own data also made the world more complicated. Having multiple devices allowed different versions of a document to be located in two places at once. A user who emails a file to himself in order to back it up has to keep track of what version it is, and how relevant it might be if he needs to restore from it. More generally, the idea of always having the “latest version” available wherever a user wants it to be available has become both paramount and largely ignored at the same time [13].

The existing state of PIM is demonstrated in an example below.

Zoe is a fairly non-technical computer user, working on a presentation using her desktop computer at work. She wants to take it home and continue working on it through the evening, so she transfers it to her portable USB drive. Using her personal laptop, she finishes the presentation at home. The next day, she’s back at work – but has forgotten her USB drive. She is over an hour away from home, but needs that file. Her two options are to drive home to get it, or recall all the changes she made the evening prior.

The most obvious use of Family might be to help Zoe see that these files are related. Still, it should be fairly clear that she already has a good grasp of this – to her, both files represent the same *concept*, and she is only interested in having the latest version available to her. Family is broadly aimed at resolving this issue – by providing these sorts of relationships to end-users and higher-level applications. However, making such grand claims about the meaning of these relationships such as claiming that two files represent the same *concept* is out of scope of our work.

We present a more futuristic vision of PIM below, corresponding to the same scenario.

Zoe is a non-technical computer user, and she is working on a presentation at work. She transfers it to her portable USB drive so she can work on it later at home. The next day, she’s back at work – and has forgotten her USB drive. She brings up the file located on her work computer, and even though she already

understands that there is a more recent version available on her home computer, she still asks her operating system to show her any outstanding or newer instances of her presentation. This interface then allows her to automatically and seamlessly retrieve and merge her content from home.

While this example may have been ‘simplified’ through the use of version control tools, these may be seen as overly technical for the average end-user. Modern operating systems do not come ‘out of the box’ as standard with this sort of functionality. They may also require explicit steps to use their features.

It is worth noting that Zoe’s initial actions in both the above examples are the same. Despite the age of the traditional file system, various research ([15; 17]) has shown that end-users are actually quite comfortable with the idioms it implies. So, while Zoe does have the benefit of resolving her problem in a new and novel way, our vision does not revolve around altering traditional user behaviour.

Existing Applications

Our introduction also defines a relationship class which describes files in unique states over time. Our motivation for this class is to reveal not only when a file was *last* changed, but also *every* time it was changed. While this information is important to create correct, time-relevant conceptual relationships – such as when a file was *copied* – it also has a simpler application to systems that already exist today.

Apple’s ‘Time Machine’ [4] is a backup tool integrated with Mac OS X Leopard that allows users to see the state of their computer’s files at regular historical intervals. However, Time Machine is limited in that it will only back up documents on a regular basis, rather than in response to changes undertaken by users. In a futuristic vision, Time Machine would back up relevant parts of a user’s hierarchy in response to when, and how often, these areas were used – a document which is saved and updated every five minutes might be more important to archive than a file which is changed every three or four days. This information is, as previously mentioned, not available within a traditional file system.

Time Machine can also not identify files with respect to their location over time. If I move **Report** to my Documents folder from my Desktop, Time Machine will simply create a new backup of **Report** in Documents and it will fail to correctly link this to the previous location that it was stored. A user then seeking an old version of **Report** on Desktop will have to remember that they moved the document at some point in the past. Ideally, Time Machine would be able to use relationship information to automatically direct users to the correct file.

1.3 Challenges

None of the aforementioned relationship classes are currently revealed through traditional file systems. Some research into PIM aims to ‘rewrite’ the way end-users perform PIM [19]. However, a core aim of Family is to capture relationships classes *without* any explicit interaction from end-users. The motivation for this aim is driven by the fact that – as introduced previously in our vision – various research suggests that users are actually comfortable with the way

the traditional file system works, and the ‘filing’ idiom it provides. And, other research ([16]) at least accepts that end-users will use the traditional idioms provided by hierarchical file systems in order to achieve their goals – even if they do not go as far as proposing that they actually like them.

Thus, we are presented with a core challenge, which reduces to determining what is occurring within a low-level file system that does not reveal this information in the first place. The file system at any given point represents the outcome of how applications have interacted with it. Thus, a natural approach to solving this problem is to analyse the actions these applications undertake while *using* the file system, and attempt to infer meaning from this.

However, this is more complicated than may be originally perceived. As mentioned previously, while users may undertake *conceptual actions* using applications – such as *save*, *rename*, *copy and paste*, etc – these same applications will actually apply the user’s conceptual meaning to the file system in terms of a multitude of lower-level system calls. Furthermore, different applications may deconstruct these conceptual actions in different ways, resulting in a situation where a single fixed approach to analysing application behaviour is insufficient.

There is also a second core challenge in supporting Family, which revolves around the idea of *unobserved changes*. While it seems difficult – but plausible – to interpret how applications deconstruct a user’s conceptual actions, and convert these to relationships, the idea of deconstructing actions from applications we cannot even formally observe appears almost impossible. The most simple example of this is if a document is sent to another user who edits it and then returns it, we would still like to identify the changes that it may have undergone. In this case, we are actually forced to observe the outcomes of behaviour rather than being able to observe the behaviour itself. However, this is both an incredibly difficult task as well as being out of the scope of this thesis. Thus, we reduce our goals: we only aim to correctly observe that a file was previously related to some file on a local computer, rather than explicitly defining *how* it is now in its current state. We can visualise this much like an ambiguous ‘cloud’ (as seen in Figure 1.1) – still providing end-users with some useful information, but of course not able to determine every point in time at which it was changed or updated, and therefore, every one of its unique states.

We discuss our algorithm and implementation – which aims to overcome these challenges – in Chapter 3.

1.4 Contributions

- The *idea* of Family

We introduce Family, a new and novel way to support users through capturing and revealing seemingly simple relationships that may already exist within a user’s *mental model*. This support is implicit, as end-users are also not required to ‘rewrite’ their PIM strategies in order for relationships between their files to be captured. Furthermore, we discuss the most relevant classes of relationships in this field, and how they map to a user view of their files through this mental model.

- Algorithm & Architecture

This thesis details how we overcome several core challenges through a concrete algorithm designed to capture information about relationships between files. This algorithm, as mentioned, is designed to observe applications in their activities over traditional file systems, rather than seeking to examine their output.

- Research Outcomes

Research and development of Family has led to various interesting conclusions about our approach to determining Family relationships.

- Implementation

We implement Family on the Mac OS X platform, using a combination of low-level code in C and high-level code in Python. This tool will also be publicly released under a GPL license, with the title *Lighthouse*. It is made available for download at <http://www.it.usyd.edu.au/~stho2688/>.

Chapter 2

Related Work

Most, if not all, modern operating systems are based around the use of traditional file systems to help users achieve their computing goals. This can be more broadly described as personal information management, or PIM. The international conference CHI defines the study of PIM as “the practice and study of the activities people perform to acquire, organize, maintain, and retrieve information for everyday use.” [1].

The traditional hierarchical file system was originally popularised through the development of UNIX [21], and to a lesser extent, Multics [10]. Daley and Neumann [12], in their work on Multics, originally laid out a specification for a general purpose file system in 1965. They discussed the idea of using symbolic addresses (i.e. filenames) to map to absolute, physical locations on random access devices such as hard disks. This paper also formalised the definition of a directory, a special file maintained by the file system which could contain a list of links to other files. This definition allowed for directories to contain other directories; Daley and Neumann simplified their definition of a hierarchical file system by describing it as “a tree of files” 2.1. Today, nearly all form of personal information management completed on any modern operating system still takes place within the scope of these ideas – originally brought into the mainstream nearly forty years ago.

Since their inception, hierarchical file systems have been both criticised [19; 7] and praised [15; 17; 8] for their limitations and benefits, respectively. Regardless, they are a core element of any modern operating systems, and form a basis for all forms of PIM.

Blunski et al. [7] introduce a vision consisting of a ‘personal dataspace’ consisting of automatically replicated emails, documents, pictures, etc. Elements of this dataspace may be explicitly shared to other users. However, all interaction, even by users external to the system, is envisioned as a seamless operation over this data which has been automatically replicated across all relevant personal devices. Blunski et al., in a large way, introduce a de-facto standard for the goals of ultimate PIM integration.

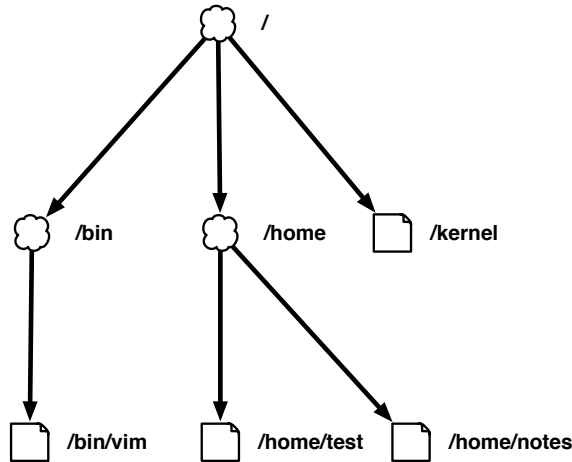


Figure 2.1: A simple file system tree, similar to the one originally outlined by Daley and Neumann [12].

2.1 Files

The traditional file system aims to provide files, ordered sequences of elements which – to an underlying file system – are essentially devoid of format [12]. No meaningful action on any of these files may be performed without, for example, the involvement of a higher-level tool such as a text editor. Various authors describe both the benefits and the difficulties with this simplistic approach. Blunschi et al. [7] coin the term “Personal Information Jungle”, reflecting the so-called jungle of currently available file formats and data processing solutions. They propose that computer users endure *physical data dependence*, relating to the fact that users need to know about devices and formats that are used to store their data in order to effectively achieve any goals related to PIM. However, it may also be argued that this unstructured approach has led to the innovation of unique and interesting file formats and storage practices, not inhibited by having to provide any higher-level meaning.

Applications will typically use files and file systems to their own advantage; storing information that they process in specific locations or in forced hierarchies. The formatless nature of files on a traditional file system provides a perfect environment for this kind of [ab]use. Furthermore, Karger and Jones [16] describe PIM as ‘fragmented’. They introduce various unification strategies, and discuss their limitations – considering that the lowest common denominator of any application is a file, and that *any* need for higher-level meaning is compromised by attempting to provide a system that supports all possible usage of this base structure.

Mahalingam, Tang and Xu [19] observe that human memory operates in terms of document content or features, rather than filenames and position in a file hierarchy (in this context, features may also represent relationships between files). They argue for new file access mechanisms, beyond hierarchies containing files of arbitrary data. This notion is similar to the approach of WinFS [20], which was a feature originally planned for Microsoft’s most recent operating system, Windows Vista, but removed from the final release. It aimed to provide

a “relational filesystem” which could provide lifelong document identifiers, as well as create relationships between these identifiers. In the initial implementation, it captured a small selection of relationships, such as ‘author’ field of an email entity being related to a unique contact card entity.

2.2 Layout

While this thesis does not focus explicitly on use patterns or behaviour of users within their own folders, they are an interesting starting point. For example, the placement of multiple files within the same folder may explicitly define – inside a user view – that they have some relationship, such as being components of a bigger task.

So, existing operating systems provide very good support for hierarchical relationship with folders. This is a valuable relationship: James et al [15] established that it serves the ‘average user’ well and Karger et al [17] show how users benefit from organising documents into a personal hierarchy. Notably, this only captures relationships if the user places related files into the same position in the hierarchy. However, a hierarchy has various limitations. Different applications may restrict or control placement of their documents within the hierarchy, or enforce their own. For instance, Mac OS X’s Mail application represents each mail message as a single file on-disk, but enforces their location and creates its own virtual hierarchy within the application.

2.3 Temporal Usage & Relationships

Current implementations of hierarchical file systems provide only limited support for users to identify temporal relationships over files. While these systems can provide end-users with the time each file was *last* changed, there is no support for retrieving information about *every* time they were changed. This information may be useful for task management or backtracking over the lifetime of one’s own documents. Additionally, having a view of changes over all files – as opposed to just a single file – may help to better represent a temporal link between changes performed on many different files during a common timeframe.

A user, Johan, may be working on a poster in some graphic design software. He opens the file and notices a typo – fixing it and instinctively saving the document. He is now interested to determine when he last worked on this document, but the ‘last changed’ field now reflects the time at which he fixed the typo. He was interested in finding the notes he was working on around the same time as the poster, based on their temporal usage, but will now have to find this file manually.

So, although files may have been changed at similar times as part of one activity, if some of them are later changed there is no way to detect the temporal link – or *relationship* – suggested by the work done on them at the same time. In this case hierarchical file systems do not support a user’s mental model of file *ancestors* as discussed initially.

2.3.1 Explicit

Cox and Josephson [11] discuss the issue of temporal usage in terms of file synchronisation. For background, Pierce and Balasubramaniam [6] present an inferred framework for describing file synchronisers – these are tools intended to manage the replication of files over multiple (potentially disconnected) traditional file systems. A conclusion that may be drawn from successful operation of these tools is the implicit relationship that the files it manages are identical to their synchronised counterparts. Cox and Josephson propose that while traditionally, synchronisation times are considered to encompass a file modification history, this is actually only a limited view of managing distinct versions of files. For example, some file synchronisation systems may perceive many changes in a single file, and a submission of this file, as a single action. However, by maintaining a distinct vector of every time a file was changed, we can better understand how all changes have been propagated to different replications prior to this point in time.

Regardless, a conclusion that may be drawn from successful operation of these tools is the implicit relationship that the files it manages are identical to their synchronised counterparts. For example, if a document titled ‘Report’ has been synchronised into two distinct hierarchical locations, these two distinct file system entities may be implicitly viewed as being an identical *version*, and exist as a single concept in a user mental model. However, even after these relationships have been identified implicitly within some software, they can not be exposed in a consistent way. No file system allows for these relationships to be defined on files directly, so these relationships only exist because higher-level tools like file synchronisers define them.

Similarly, version control systems [9; 18], widely used by programmers, can assist in identifying relationships between existing files and their previous versions. As well as this, these forms of tools can provide temporal information about ‘changesets’, or blocks of work – potentially completed on multiple files – completed at specific points in time. In a similar vein to file synchronisers, these changes must be explicitly *checked in* to the system, and will treat all work completed since the last check in as a single block, ignoring specific temporal changes.

However, users must interact explicitly with both traditional version control tools and file synchronisers. For example, a user must explicitly stop and *check in* their changes. Mainstream version control tools are not able to uniquely identify user work and implicitly track this over time.

2.3.2 Implicit

Conversely, there has been some work towards our vision of implicitly capturing and presenting relationships. Soules and Ganger [22] created a system capable of making use of information about file system access to identify temporal relationships between files. Essentially, they aimed to group sets of files based on the time at which those files were used. This grouping was then used to enhance search results for end-users. They reported that this approach improved search precision from 17% to 28%. They also made use of the ‘relationship’ idiom, creating a relation graph based on file system usage. However, their relationships did not include the concept of *versions*, or our notion of *ancestor*. For them, every file was treated as a single, most current version.

Chapter 3

Approach

This chapter describes the the algorithm behind Family. The primary contribution of this section is the description of the algorithm, and architecture, used within Family to capture conceptual relationships. However, we will first introduce the relevant terminology, as these terms are unfortunately fairly overloaded.

We introduce three key terms: *document*, *file* and *node*. These are important as they bridge the gap between a user's *mental model* and the extremely raw file system view describing the unique state of every file.

- *Mental model*

Broadly, a user's *mental model* describes the way in which a user thinks about each file and its relationships with other files. For example, a user might consider several documents in various folders to be related to each other. One of these files may have originally come from an email attachment, and a user's mental model might include this relationship. Another file, **Notes new**, may be a copy of an older file, **Notes**, which has been changed. However, modern operating systems do not typically attempt to support this model. They do not understand or attempt to support conceptual meaning or relationships such as 'new versions'; to them, every document is just a *file* – an independent sequence of bytes with a meaningless name.

- *Document*

The notion of *document* is part of a user's mental model. It will have a known name, such as **Report 10-02**; a known location, such as a user's desktop or on a USB stick; and typically also contain content that may be viewed or changed by an application able to open documents of its type. And, regardless of how a document is actually technically contained on a disk, it will typically be seen by a user as a single entity that remains constant between its creation and deletion. If the document is copied, duplicated or created by saving another document with a new name, a user is able to update their own internal mental model to capture the way in which these new files are related.

- *File*

The term *file* tends to be overused, but from the view of the operating system, a file can be simply defined as a filename and a sequence of bytes. From the higher-level view of an application, a *file* may be interpreted any way they like – a PDF viewer will interpret a PDF in order to render it, while a tool to count its pages will ignore much of the content and focus on a few specific facets.

From the view of an end-user, the terms *file* and *document* tend to represent the same idea. That being said, we have throughout this thesis used both terms fairly interchangeably in regards to the way users view their file system. However, while a *document* can be understood and arranged by a user into relevant mental models or user views, a *file* does not enjoy the same privilege – it only exists in a technical context. Additionally, a *document* may be comprised of a package folder which actually encompasses multiple files.

Modern operating systems and file systems also support the notion of extended attributes, which may be applied to any specific *file*. These are arbitrary key/value pairs that may enhance a file or provide out of band data. Additionally, one of the fundamental properties of files is that they are mutable, and only have a current state. They are created and updated over their lifetime backed by a single *catalog identifier*, a unique identifier that has explicit relevance to each specific mounted file system. Users may save new versions of files regularly, causing the overall file concept to be updated, and information about the previous state to be lost.

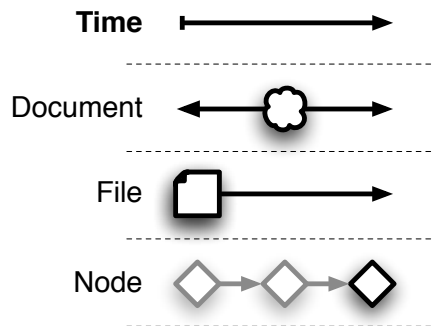


Figure 3.1: Demonstrating the leap between our three core concepts.

- *Node*

Family embraces the term *node* to describe an immutable state of a file, where *state* describes a file’s name and content for a particular period in time. A file must always map to at least one node, as its content always represents one state. However, a file may also map to multiple nodes, if it has previously existed in other states. This is the case if, for instance, a file undergoes changes – rather than being created and deleted but never modified, in which case the file will continue only to map to a single node. This distinction is demonstrated explicitly in Figure 3.1.

Technically, each *node* has a small set of properties which allow it to be mapped to

a real-world *file*. These include properties such as file size, local *catalog identifier*, along with the file’s creation and last modified times. Each *node* also contains a set of evidence, beginning empty, which may grow throughout its lifetime. These facets of evidence either provide information about the node itself (such as when it was last opened), or they provide information about how this node relates to others (such as if it is a *descendant* of a previous node). Both these important features of nodes – properties and evidence – are heavily used within our algorithm, introduced below.

We show the jump between our three core definitions – of *document*, *file* and *node* – in Figure 3.1.

3.1 High-level Formalism

Firstly, we describe various formal approaches for the form of data we aim to capture. Our evaluation, to be completed in Chapter 5, is assessed on the basis of our implementation’s correctness versus these approaches.

- $Anc(X)$ – or ancestors – *nodes* the user acted on to create X
- $Dec(X)$ – or descendants – *nodes* of which X is an ancestor
- $Family(X)$ – a combination of both $Anc(X)$ and $Dec(Anc(X))$.

These formalisms are key to Family, as they describe in specific terms our ultimate goals. While we have originally described relationships in terms of classes, these ultimately reduce to either providing information about *ancestors*, or information about *descendants*. It is also worth noting that the two base cases, $Anc(X)$ and $Dec(X)$, are simply the inverse of each other.

We can list a short example of *ancestors* of an arbitrary node, A . Content may have been copied from B into A – therefore, B fits our definition of ancestor. However, A may have also been originally created by copying C . Additionally, C may have been originally created by simply making changes from D ! Thus all these nodes – B , C , and D – would fit the definition of $Anc(A)$.

Conversely, we can show the *descendants* of node B as much the inverse. As content was copied from B to A , it fits the definition of a descendant. Additionally, if B was at one point copied to F , it also gains the label of descendant. So A and F both fit $Dec(B)$.

3.2 Conceptual Models

The above formalisms are constructed through *node* relationships that fit various definitions within a user’s traditional mental model. The below conceptual models are the most important set of these relationships, and they explicitly define how the nodes they involve should be defined in terms of ancestors and descendants.

We have identified four important cases described below.

- Document copy or ‘save as’

For example, a user might use a command-line tool to copy one file (e.g. **Notes**) to another, brand new file (e.g. **Abstract**). They may alternatively open **Notes** with a text editor and use the standard ‘save as’ functionality to create the new file, **Abstract**. By mapping this user functionality to abstract nodes, relationships may be drawn that would allow users to backtrack from the newer file, in this case **Abstract**, to its original source – **Notes**. Formally, the *ancestors* of the node behind **Abstract** now include the node behind **Notes**. Similarly, the *descendants* of **Notes** now include **Abstract**.

- Document changes (edits) over time

As a user saves his or her work in-place over time, each save is identified by the creation of a brand new *node* representing the current state. This provides support for identifying revisions and historical relationships. For instance, if a document is copied at 10:15 to create a new version located elsewhere, it is important to appreciate that this copy was made from a specific immutable state of a file (i.e. a *node*), rather than from its entirety as a mutable concept over time.

Again, the way this conceptual model fits within our definitions of *ancestors* and *descendants* is fairly clear-cut. The older state is now a direct *ancestor* of the current state, and the current state is now a direct *descendant* of the older state.

- Related work on documents, including copy and paste from one document to another

It is important to identify more general relationships between two documents. Documents can be related on a task level – for instance, two documents may be open at the same time as they are both related to achieving the same user goal. Text and content may even be copied back and forth from both documents to each other. By identifying that these documents were in use at the same time, we support the notion that users can backtrack from their current state and determine documents which are not strictly related in the sense of *family*, but are related in other ways and may still be useful to support the mental model.

More specifically, this relationship is again very back-and-forth. If content is copied from one document to another, then the original document must now be perceived as an *ancestor*, while the target document must be considered as a *descendant*.

- Documents sent and returned outside the machine

The previous three examples all refer to specific cases that may be observed on a controlled, local system. However, Family extends to work performed on multiple machines or by multiple users. Put simply, a user may send a file to a colleague, who may make changes to the file and return it. A user’s mental model might support the idea that this file has a similar name, and has a later modification date, so it is likely that this new file is somewhat related to the original. We would like to identify this relationship semi-formally, with an understanding that while the file is ‘out of control’, its state is generally ambiguous (represented by a cloud in Figure ??).

While this conceptual model is less specific, it still fits into our definitions of *descendants* and *ancestors*. This is largely because despite the uncertainty of the actions a file may have undergone outside the system, there is still a temporal connection between the

original file and the newer, observed file. In a broad sense, we see the change as one ‘leap’, as opposed to many, but we can still see the returned file as an *ancestor* of the last state it was observed in. Furthermore, like our previous examples, the opposite point of view considering *descendants* is also true.

3.3 Low-level analysis

The algorithm of Family, in an overall sense, involves observing system calls in order to determine relationships between *nodes*. These relationships aim to map to the conceptual actions users have undertaken, such as *save*, *copy* and *duplicate*, which then map to the base cases of *ancestors* and *descendants*.

System calls

We perform this algorithm for every single system call we observe. These system calls may be slightly different depending on the operating system focused on, but they fit into several simple broad categories. Many, for instance, are merely useful to suggest that a file should be analysed for changes: even though they do not explicitly define that a change has actually occurred. They fall into the following three categories.

- Requests; e.g. `open()`

These include events such as the low-level action of opening a file, including authorisations for specific actions. For instance, a PDF reader may open a file requesting only to read it. However, a text editor may open a file with both read and write authorisation. These authorisations do not imply that an application will modify a file, just that they have *permission* to do so.

- Changes; e.g. `unlink()`, `creat()`, `write()`, and varieties of `open()`

Single-file changes come in a variety of forms. In most cases, these will be unique and will be recognised by Family in various different ways. Broadly, they fall into several sub-categories.

- File creation

The creation of files is a key action that may be undertaken by applications. In some cases, this event may be split into the creation of a file versus the creation of a folder or directory. There may also be a special case if a *link* is being created to a previous file, rather than dealing with a brand new file.

- File deletion

Deleting files is another key action that may be undertaken by applications. Again, special considerations must be made where the file is linked to from multiple locations. Files may also be deleted as the result of rename operations, discussed below.

- Changes in content

Depending on how a file system provides these events, they may be provided in a single notification after-the-fact or as intermediate notifications on specific reads. Regardless, these events are typically coalesced as an application may perform hundreds or thousands of individual changes in order to perform an overall update.

- Changes to extended attributes

A file may have its owner changed, or extended attributes changed. These changes are less relevant to Family, as most applications do not perform substantial activity using these operations. In fact, Family actually uses extended attributes to achieve its goals, therefore creating an issue with loopback if these changes (which Family performs itself) are observed by itself again.

- Movement or structural change; e.g. `rename()`, `exchange()`

As well as changes to files in terms of their creation, deletion and content, we are also concerned with their movement within a system. If a file is renamed to a new, unique name, this is perceived as a new unique state. There is also a special case where an operating system attempts to rename *over* an existing file. This is interesting – as regardless of its underlying meaning, it implies that the existing file is to be removed.

Mac OS X also provides for a unique system call: `exchange()`. It allows for the content of two files to be atomically ‘flipped’, therefore invalidating the current state of both. We support this system call within the architecture of Family, but again, it is only provided on Mac OS.

Since our algorithm (defined in 3.4) analyses every single system call, these system relevant as they allow the algorithm to perform differently depending on what category they fall into. The quirks of each event are also vital to understand in order to correctly process them.

Node identification

From a functional point of view, Family attempts to identify and ‘tag’ all files currently in use by a user. Tagging a file allows it to be uniquely identified over time, and in fact, each tag maps directly to a node. As files are mutable over time, Family seeks to observe any change in its properties which may imply an overall change in state. It performs this by attempting to assert that the newly changed set of properties still matches the properties tied to the given ‘tag’ and its mapped node. If the assertion fails, a new node is created – marking the old node as destroyed – and evidence is added to give reason to the way the new node is related to the old node. In the simplest case, this evidence simply lists that the new node is directly derived from the previous node. The properties that may be examined include a unique identifier for the local disk, the file’s unique identifier (or ‘catalog ID’) on disk, a unique identifier for the file’s containing folder, its size, its creation date, and its last modification date. This set also includes the file’s name.

For example, a pre-existing node may be copied (using a file manager like Mac OS X’s Finder) to a new node, creating a brand new file. In this case, the pre-existing node is never changed, but Family will detect that the new file was created from copying the pre-existing node, and

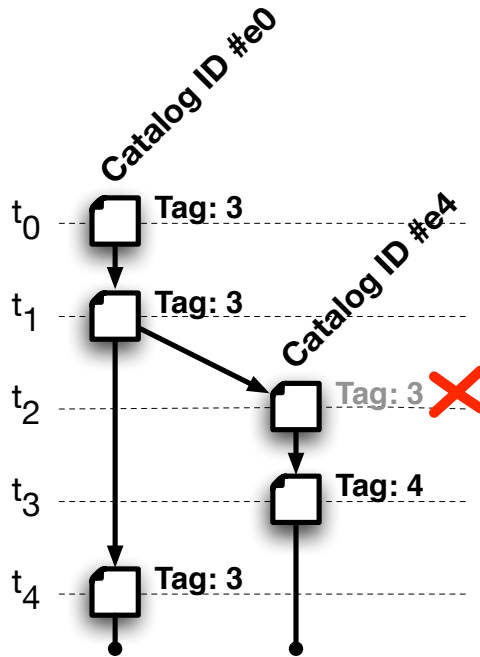


Figure 3.2: Identifying a copied file

create a new node with evidence to suggest as such. Figure 3.2 demonstrates this example making reference to the ‘tag’ idea described above. A file is first created at t_0 , located within the local disk with catalog ID #e0. Family then assigns it the next available tag, 3. It is edited at t_1 , retaining the same tag. At t_2 , the original file is copied along with its tag to a new location, catalog ID #e4. As the tag 3 does not correspond to its known catalog ID, Family assigns the file a brand new tag, 4, to this new file. This change is undertaken at t_3 . As well as updating this tag, Family may now (at a higher level) create evidence linking these two unique states of both files. The original file remains unchanged.

3.4 Algorithm

Our algorithm is listed below. It interprets each file system event (each falling into one of the categories and/or sub-categories above) individually and attempts to correctly construct nodes and derive relationships through an ongoing supply of these events.

For each observed event, we:

- Attempt to map the *focus* of this event (a *file*) to a *node* that is already known
 - i. Look up the node that this focus file maps to, through the node identifier stored as an extended attribute on this focus file itself
 - ii. If no identifier was found, or its corresponding node was not found, we have not mapped this file

- iii. Otherwise, the ‘correctness’ of this mapping is determined by whether our known properties (e.g. file size, modified time) match the properties found on the focus file
- Use this mapping to determine the most correct *node* to work on
 - If the focus file already maps to a node, and this mapping is ‘correct’
 - i. Do not update or change existing nodes
 - If the mapping is ‘incorrect’:
 - i. Create new node – mark the focus file as corresponding to this node through an extended attribute
 - ii. Link this new node to the previously mapped node through new evidence

This evidence will either imply that the file has been *copied* (a relationship between files), or simply *updated* (a temporal relationship representing unique states over time).
 - If the focus file is not mapped at all:
 - i. Create a new, unrelated node with no evidence – again, mark the focus file as this node through an attribute
- Apply evidence relevant to the event type

While many events are essentially advisory, and do not imply any specific evidence we should create, there are several key types of events that the algorithm behind Family will look for.

- i. The deletion event, whether explicit or caused by a rename

Evidence will be submitted to the target node marking it as removed. In the case of being overwritten, evidence will be submitted to the file being renamed indicating that at this point in time, it was responsible for the demise of the now deleted node.

Chapter 4

Implementation

Family is primarily implemented as a user-space Python tool, named Lighthouse, which runs on Mac OS X Leopard [5]. It sources various low-level information from its current environment, such as system calls, and applies the algorithm described in Chapter 3 in order to identify and represent relationships between unique file states.

4.1 Kernel Extension

Lighthouse relies heavily on a Mac OS X kernel extension, named `deveye` and written in C, developed by the author of this thesis specifically for this task. Further information about the public interface to this extension is available in Appendix A. It interfaces with `fsevents` [2] to provide real-time information about all file changes on a Mac OS X system, along with a custom-built interface to `KAuth` [3] which provides real-time information about current file usage (including notification of file reads).

An additional feature of `deveye` is that it allows user-relevant access to this information, as normally `fsevents` and `KAuth` are system-level devices only accessible by the superuser. This information is revealed through a public API, implemented through `ioctl()`, that may be called by all users on the system to find events relevant to that specific user only – without any special privileges. So, while the kernel driver itself must still be installed with superuser privileges, once it has been installed any application run by a unprivileged user may use its output – including the program written for this thesis which runs completely in user-space.

Again, the installation and usage of this extension is briefly detailed in Appendix A. It is a new and useful software tool which has considerable value for exploring new ways to support user information. For example, it has served as a foundation for another Honours thesis project in 2008 which aims to support users in taking an activity-based view of their file system.

4.2 Lighthouse

Firstly, the Python implementation of Lighthouse uses a small C to Python bridge known as *eyedropper*. This simply provides information from the kernel extension, accessible through `ioctl()`, in a Python-like form for rapid development. Lighthouse is almost a direct implementation of the algorithm described in Chapter 3. It functions by parsing each system call, determining the relevant focus *file* and mapping this to a relevant *node* representing that file's unique state. It also embraces the three categories of events to determine the most appropriate relationships to create. These relationships are then stored using `SQLObject` [14] inside a small file-based database, available for any further tools to examine and parse.

Most notably, this implementation does not yet correctly support the conceptual model of 'copy and paste', or the operation of 'save as'. This is due to our limited time during our research into the area, and the vast many ways applications attempt to achieve these simple goals. Much research into enhancing PIM focuses on specific applications, such as popular modern office suites. However, we have been hesitant to analyse specific applications as it – in the long term – severely limits the focus of any research.

4.3 Testing Framework

We have also developed a Python tool which uses the data produced by Lighthouse in a resolution process designed to provide broad information about the *ancestors* and *descendants* of any specified file. This tool is a simple, publicly available user interface to what Lighthouse generates, and it may also be used as the basis for our testing harness which will be described in the next chapter.

In terms of our formalisms, this tool attempts to map a given *file* to a current, active target *node*. This *node* represents the current state of the given *file*. It then resolves the output of $Anc(X)$, $Dec(X)$ and $Family(X)$ by inspecting the available evidence attached this *node*, and the further *nodes* this evidence leads to, etc. The *nodes* found through this process are then filtered to include only those which are the most active, and which currently exist – that is, *nodes* which map to *files* that a user can currently access.

This mapping example is visualised in Figure 4.1. The *files* on the left of the figure indicate items that the user has current access to. These are the user's files: **Notes**, **Report** and **Abstract**, shown as clipped squares. For the example, the file **Notes** has three *nodes* over time, the first version at t_0 , the next at t_1 , and the most recently changed version of the node was created at t_5 . This latest version corresponds to the user-accessible file, the last time in the figure.

Similarly, the file **Report** has three *nodes* over time – at t_2 , t_3 and t_6 . **Abstract** also has unique versions at t_4 , t_6 and t_8 .

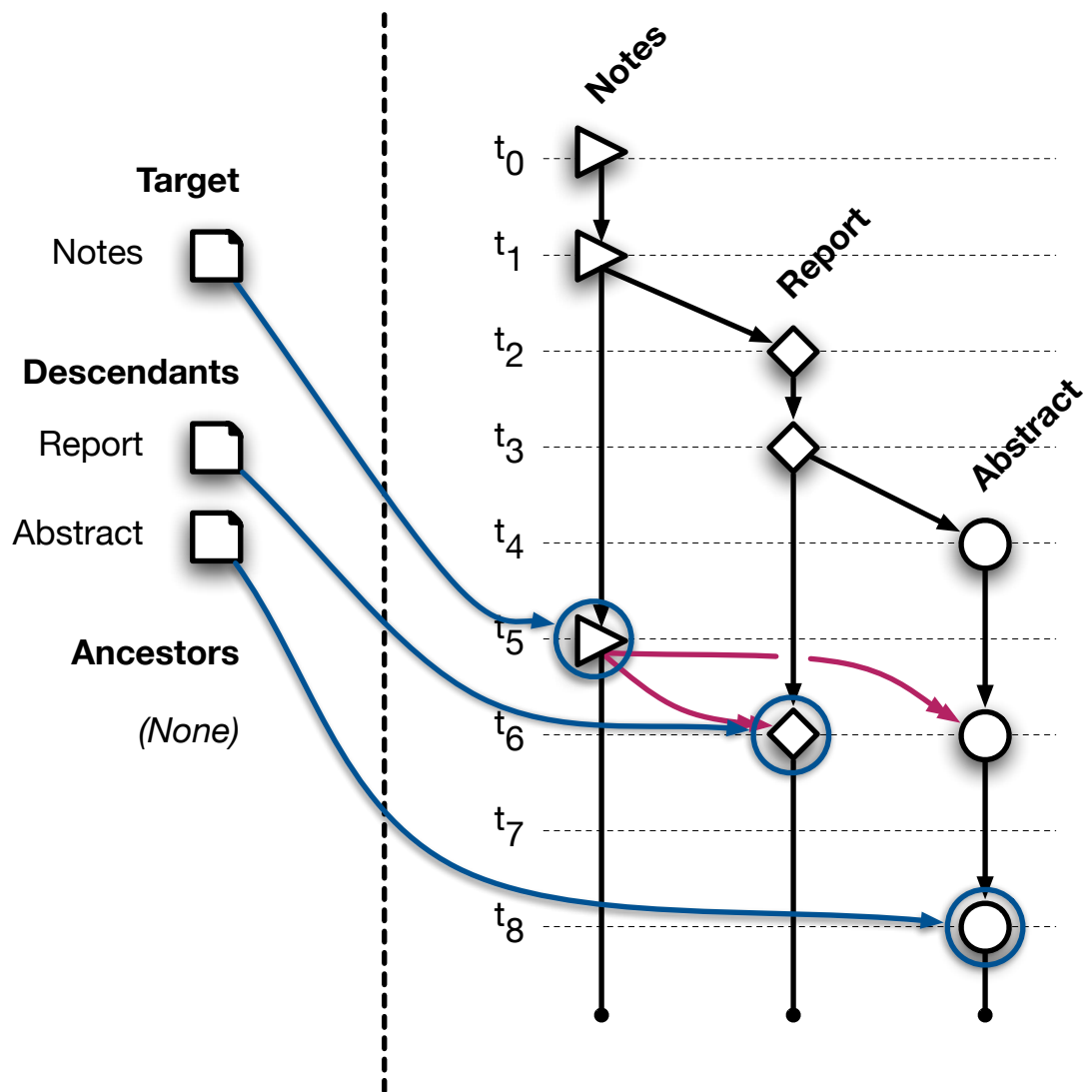
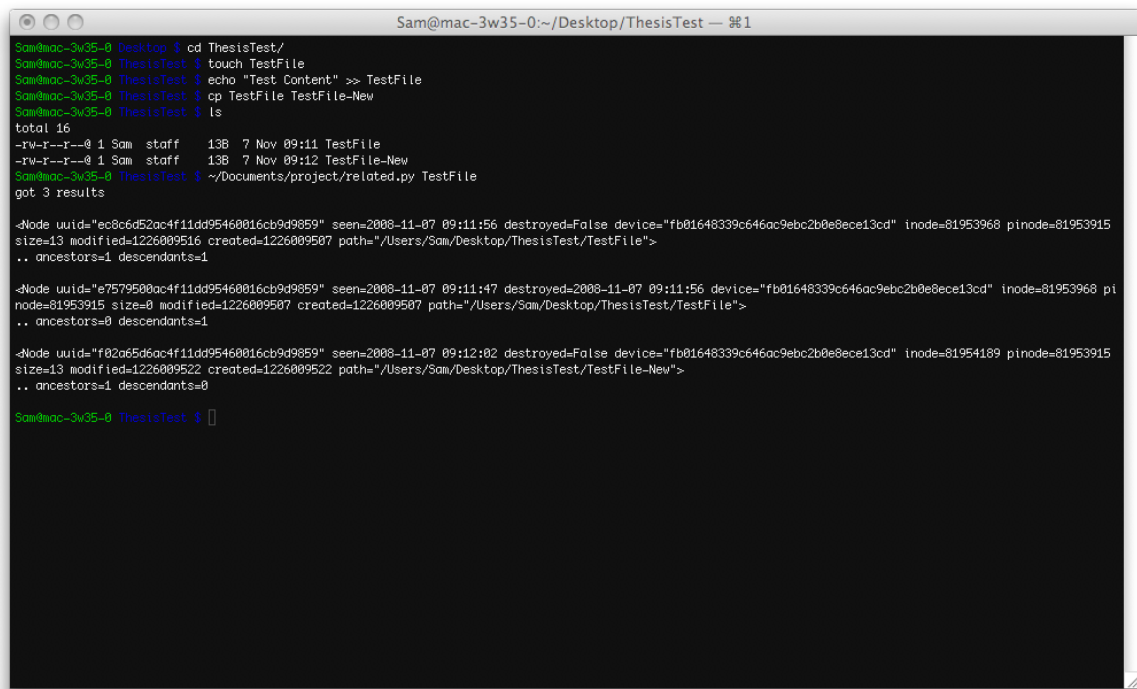


Figure 4.1: Mapping files to nodes, as well as finding their ancestors and descendants.

4.4 Walkthrough

In Figure 4.2, we demonstrate a simple example of how the testing framework is able to reveal a simple file copy, through a small walkthrough in the Mac OS X [5] Terminal. In this example, **TestFile** is created (creating a brand new node) and then a small amount of content is appended to it (creating an updated state of **TestFile**, and therefore another node). The file is then copied to **TestFile-New**. We then run our testing framework on one of the files within the scenario, revealing every node, including the original state of **TestFile** – which is now marked as destroyed, as it has been replaced with a more current state.



```
Sam@mac-3w35-0:~/Desktop/ThesisTest — 1
Sam@mac-3w35-0 Desktop $ cd ThesisTest/
Sam@mac-3w35-0 ThesisTest $ touch TestFile
Sam@mac-3w35-0 ThesisTest $ echo "Test Content" >> TestFile
Sam@mac-3w35-0 ThesisTest $ cp TestFile TestFile-New
Sam@mac-3w35-0 ThesisTest $ ls
total 16
-rw-r--r--@ 1 Sam  staff  13B  7 Nov 09:11 TestFile
-rw-r--r--@ 1 Sam  staff  13B  7 Nov 09:12 TestFile-New
Sam@mac-3w35-0 ThesisTest $ ~/Documents/project/related.py TestFile
got 3 results

-Node uuid="ec8c6d52ac4f11dd95460816cb9d9859" seen=2008-11-07 09:11:56 destroyed=False device="fb01648339c646ac9ebc2b0e8ece13cd" inode=81953968 pinode=81953915
size=13 modified=1226009516 created=1226009507 path="/Users/Sam/Desktop/ThesisTest/TestFile">
.. ancestors=1 descendants=1

-Node uuid="e7579500ac4f11dd95460816cb9d9859" seen=2008-11-07 09:11:47 destroyed=2008-11-07 09:11:56 device="fb01648339c646ac9ebc2b0e8ece13cd" inode=81953968 pi
node=81953915 size=0 modified=1226009507 created=1226009507 path="/Users/Sam/Desktop/ThesisTest/TestFile">
.. ancestors=0 descendants=1

-Node uuid="f02ac65d6ac4f11dd95460816cb9d9859" seen=2008-11-07 09:12:02 destroyed=False device="fb01648339c646ac9ebc2b0e8ece13cd" inode=81954109 pinode=81953915
size=13 modified=1226009522 created=1226009522 path="/Users/Sam/Desktop/ThesisTest/TestFile-New">
.. ancestors=1 descendants=0

Sam@mac-3w35-0 ThesisTest $
```

Figure 4.2: Walkthrough of the testing framework

Chapter 5

Evaluation

We evaluate Family by performing a series of synthetic tasks and comparing their output to the correct output, in terms of *ancestors* and *descendants*, and hence the overall notion of Family. By this process, we evaluate the correctness of our algorithm and implementation.

We also evaluate our implementation of Family in terms of system runtime performance, in two key areas. Firstly, we determine the ongoing baseload cost of capturing and processing events in real-time. Secondly, we analyse the time it takes from when a user performs a conceptual action to when this action is correctly mapped in a way that it can be later retrieved. These need to be fast since they constitute an ongoing background cost for this approach.

We do not assess the cost in ‘looking up’ family relationships – while lookups may be relatively costly, they are provided purely as a benefit to end-users. Thus, users only incur cost when they are interactively attempting to find more information about their files, a cost they may well expect. However, manual testing as shown this cost to be on average less than one second for the vast majority of look ups.

5.1 Correctness

To determine the correctness of Family, we designed and implemented a testing harness. We designed a series of tasks, utilising different application and operating circumstances, and define each as a concrete test. These tests were then be performed in terms of the previously defined formalisms, $Anc(X)$, $Dec(X)$ and $Family(X)$. More explicitly, we aimed to assess whether for each node, z , of our test set, that $Anc(z)$, $Dec(z)$ and $Family(z)$ return nodes that fit a set of proposed, expected output.

Single case tests

We designed a series of tests that are pure tests of the key elements of our conceptual model, each based on a minimal set of files undergoing a single unique action one or many times.

While some of these tests may operate on similar file names, we conduct the tests with a ‘clean slate’ at the beginning of each.

1. Copy using Mac OS X Finder [5]; **Notes** to **Report**

This is a fairly simple case of duplicating a file. Much of the motivation and development behind Family has come from the understanding of this core relationship. As in our ‘copy’ conceptual model, this action defines the newer file as a descendant of the older, and the older file as an ancestor of the newer – the older file was used to create the newer file.

2. Copy folder using Mac OS X Finder; **NETS.key** to **Dup.key**

In this case, all files within one folder should correspond correctly to their counterparts in the other. **Dup.key** files should match their counterparts in **NETS Presentation.key** as their ancestors, and vice versa. This is an extended version of the first test, designed to evaluate whether location proximity and temporal proximity (as all files will be copied at a similar time) will not affect the correct outcome.

3. Copy and paste using TextEdit [5]; **Notes** to **Report**

This test should result in identifying **Report** as now descendant from **Notes**, while **Notes** is now a new ancestor of **Report**. In many ways, this is similar to our base case of copying files, except in terms of specific content rather than a file concept as a whole. This is important as it represents one important and open way that files are created and manipulated within applications. For instance, users must be able to create a document either by copying it in an operating system interface (as per our first test), but they may be equally as happy to copy an entire file’s content into a new window and save that. From an end-user’s point of view, they may see no conceptual difference in what they are aiming to achieve, even though this test may be much more complex than a simple file copy operation.

4. Save as using SubEthaEdit [23]; **core.c** to **core-v2.c**

This tests another conceptual action which – while implemented in a more complex way – is potentially no different from the case where an end user than performs a simple file copy. Regardless, it has slightly more breadth than a simple copy. A user may make changes before saving a new copy of the previous file.

5. Save using SubEthaEdit; **Report**

We will expect to see many states of **Report** over time, corresponding to every time the file is saved to disk. The ability to correctly identify this relationship forms the basis for all temporal information between nodes in Family. It is therefore extremely useful to confirm that this relationship is created successfully. This matches the second conceptual model, listing changes over time.

6. Save using vim; **core.c**

As the command-line tool **vim** interacts with most file systems in a very unique way, it is a worthy candidate to test. However, for all other purposes, **vim** will function in the same way as the above test on SubEthaEdit.

7. Backup using `cp`; **Archive.tar** to multiple locations

Using the command-line utility `cp`, we will determine that all backups from one source file correctly include the source file as one of their ancestors. While this test is simply an expanded version of the first copy test, this is important as it ensures that the broad notion of Family is maintained. In this case, our original file must have a number of descendants all directly attached to its node.

8. External changes; **Birdie.jpeg**

We describe this simple external example as an expanded form of saving files over time. This test involves a file – located locally on a USB stick – being taken from the system and having changes made to it. In the case of **Birdie.jpeg**, a photo, we might expect the image to be touched up, or cropped etc. We want to be able to determine that the file has undergone changes, and represent this in a temporal fashion, but we are not concerned with every specific change.

We summarise these single case test results in the first three columns of Table 5.1. This table shows each test, along with the application that was used to perform it. It then describes every file that should, or otherwise is, related to this task – including the file that could be either described as *source* or *destination*. More technically, we also show each unique state of each file, with the ‘current’ and available version indicated in bold. For instance, in the case of saving an existing file, we represent the previous unique states of both **Report** and **core.c** as unbolded listings of the same file name.

Most importantly, Table 5.1 shows the both the expected and actual count of *ancestors* and *descendants* that are attached to each node. While these, of course, link to other nodes, these are correct within the bounds of each test – that is to say, the descendants and ancestors listed for each node only reference other nodes within the same test.

Results

In terms of actual results, all but the 3rd and 4th single case tests have been successful. Our expected set of ancestors and descendants are correct for all cases except that of copy and paste, and save as. While this relationship would be important for the complete success of the Family algorithm, as discussed within Chapter 4, the implementation does not currently extend to supporting this kind of conceptual model.

5.2 Performance

In terms of performance, Family runs as a subtle and ongoing background process, only undertaking activity as a result of other user actions. In this sense, it is difficult to perform a performance analysis - Family involves several processes and the costs are incurred as an effect of each user action, as well as other system background actions. We have allowed Family, as well as the low-level kernel extension, *deveye*, to run in parallel with our normal activity of an iMac 2.16ghz Core 2 Duo, with 2gb RAM and there has been no noticeable degradation in the responsiveness of the system under a diverse range of desktop activities.

Test	Application	File	Expected		Results		Tot.
			<i>Anc.</i>	<i>Dec.</i>	<i>Anc.</i>	<i>Dec.</i>	
Copy	Finder	Notes	0	1	0	1	2/2
		Report	1	0	1	0	
Copy Folder	Finder	NETS.key	0	1	0	1	6/6
		NETS.key/PkgInfo	0	1	0	1	
		NETS.key/Index.gz	0	1	0	1	
		Dup.key	1	0	1	0	
		Dup.key/PkgInfo	1	0	1	0	
		Dup.key/Index.gz	1	0	1	0	
Copy & Paste	TextEdit	Notes	0	1	0	0	0/2
		Report	1	0	0	0	
Save As	SubEthaEdit	core.c	0	1	0	0	0/2
		core-v2.c	1	0	0	0	
Save Existing File	SubEthaEdit	Report	0	1	0	1	4/4
		Report	1	1	1	1	
		Report	1	0	1	0	
Save Existing File	vim	core.c	0	1	0	1	4/4
		core.c	1	1	1	1	
		core.c	1	0	1	0	
Multiple Copies	cp	Archive.tar	0	5	0	5	10/10
		<i>Others (x5)</i>	1	0	1	0	
External Changes	<i>n/a</i>	Birdie.jpeg	0	1	0	1	2/2
		Birdie.jpeg	1	0	1	0	

Figure 5.1: Single case test results

Baseload cost of Family is shown in Figure 5.2. Broadly, this shows no substantial difference in three different environments. The first environment represents the base load environment, without any component of Family. The second environment includes the kernel extension, `deveye`, along with its support utility. The last environment also includes the user-space tool Lighthouse. The only notable feature of any environment is the tendency for the complete Lighthouse tool to present a greater standard deviation in its results.

This test suite was contributed by Greg Darke.

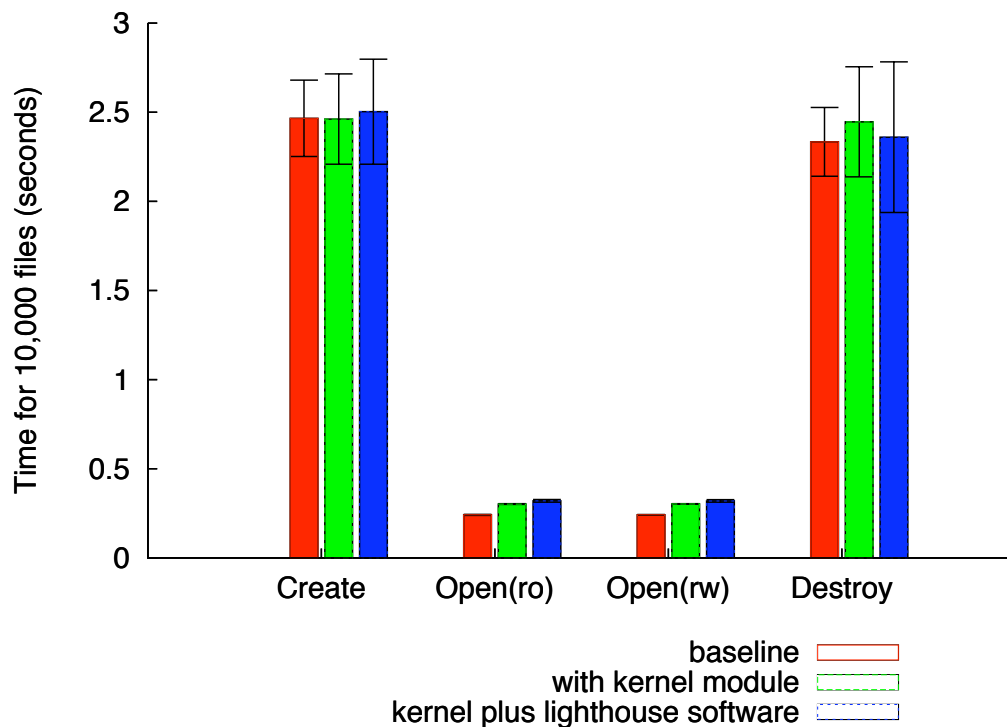


Figure 5.2: Timing for traditional file system actions

We also performed short analysis of the time taken, reported by Lighthouse, to process sets of various actions. We show this in Figure 5.3. This is the total time it takes Lighthouse to correctly parse each system call correctly – this does not count any time that it sits idle waiting for further events. As these tests were performed on a live system, the total number of system calls parsed during the time the formal test was being conducted may be much higher than that of the actions taken, due to other background processes. Despite the relative high speed of processing single events, approximately 60% of all work involved in Lighthouse is taken up by database usage, something that could very easily be optimised in the future.

Action Type	Action Count	Internal Events	All Overhead		Database Only	
			Time	Per Action	Time	Per Action
Copy	500	1284	8.5318s	0.0171s	5.3114s	0.0106s

Figure 5.3: Time taken for correct mapping

Chapter 6

Conclusion

6.1 Contributions

This thesis has explored and contributed a new and novel way to help users find files that are related to each other because of their common ancestry or usage patterns. This is valuable as these relationships may be important to help users find the right file for a particular task, or when an activity calls for key related files that were used to create an existing file.

The key contributions of the thesis include the introduction of a new notion, which aims to capture relationships that exist in a user's mental model but are totally ignored in current operating systems. We have established a formal description of this user notion of a Family, in terms of the concepts, *ancestors*, *descendants* and unique states – or *nodes* – of all files. We have also demonstrated that it is practical to correctly capture many forms of conceptual models, and classify these relationships in terms of *ancestors* and *descendants*. We have also contributed an implementation and test cases for this implementation. Lighthouse is a functional and usable tool that provides a local database, able to be queried. And it is implemented – as originally intended – as an implicit and transparent background tool. The performance of the system behaviour while Lighthouse is running is also negligible versus a base case. Thus it captures data with no ongoing cost to end users, providing information when users or upstream applications request it.

6.2 Future Work

This work provides foundations for a whole new field of research into file access interfaces. This will include user studies, determining the usefulness of having this *ancestor* and *descendant* information accessible in a easy-to-use and transparent fashion.

One broad conclusion that we have reached is that operating systems in general need to have a better conceptual understanding of user *intent*. This can be described as a ‘bottom-up’ approach to problem solving. Many typical actions on a machine are so implicit, so everyday, that it seems extreme to say that the internals of an operating system should be able to

have a handle on how or *why* they are performing a certain action. This has implications not only for attempting to observe a modern system, but also for security and other areas.

For example, the distinction between how a user loads a program could be extremely useful in providing relevant warnings or alerts to end-users. A certain level of security rating could be applied to a program if a user physically uses a mouse click to start it, or explicitly types it from the command line. However, the same level of security might not be guaranteed if the same program was started by a potentially malicious background process.

A similar analogy can be applied to our work. We aim to analyse events ‘after-the-fact’, or from a ‘top-down’ approach – which, after this thesis, should be understood as a substantial effort. If applications could inform the operating system about the simple, conceptual actions they intend to take, then this higher-level conceptual information could be mined with ease. Instead, these same applications may now embrace quirks and intricacies in order to achieve their goals, such as ‘guaranteed’ atomic file saves. However, the real issue is that nearly every application embraces the *same* quirks. If operating systems were updated to provide ‘quirk-free’ interfaces (even if these interfaces were eventually implemented in a quirky way), then information about calls to these interfaces could be analysed – a more ‘top-down’ approach.

6.3 Conclusions

Existing file systems fail to support users in identifying important relationships between their own documents that they may hold within a personal mental model. It is not straightforward to create a system that can recognise these relationships. We have explored the algorithm and architecture of Family, as well as a real-world implementation in Lighthouse, which aims to work towards this goal.

Bibliography

- [1] Personal Information Management: PIM 2008 Workshop. <http://pim2008.ethz.ch/>, 2008.
- [2] Apple, Inc. File System Events Programming Guide: Introduction. http://developer.apple.com/documentation/Darwin/Conceptual/FSEvents_ProgGuide/Introduction/chapter_2_section_1.html.
- [3] Apple, Inc. Technical Note TN2127: Kernel Authorization. <http://developer.apple.com/technotes/tn2005/tn2127.html>.
- [4] Apple, Inc. Mac OS X Leopard - Features - Time Machine. <http://www.apple.com/macosx/features/timemachine.html>, 2007.
- [5] Apple, Inc. Apple – Mac OS X Leopard. <http://www.apple.com/macosx>, 2008.
- [6] S. Balasubramaniam and B.C. Pierce. What is a file synchronizer? *Proceedings of the 4th annual ACM/IEEE international conference on Mobile computing and networking*, pages 98–108, 1998.
- [7] L. Blunschi, J.P. Dittrich, O.R. Girard, S.K. Karakashian and M.A.V. Salles. A dataspace odyssey: The iMeMex personal dataspace management system. *Proceedings of the Conference on Innovative Data Systems Research (CIDR)*, 2007.
- [8] R. Boardman and M.A. Sasse. “Stuff goes into the computer and doesn’t come out”: a cross-tool study of personal information management. *Proceedings of the SIGCHI conference on Human factors in computing systems*, pages 583–590, 2004.
- [9] CollabNet. SCM Category Home. <http://scm.tigris.org/>, visited 2008/09/28.
- [10] Fernando J. Corbató and V. A. Vyssotsky. Introduction and Overview of the Multics System. *IEEE Ann. Hist. Comput.*, Volume 14, Number 2, pages 12–13, 1992.
- [11] R. Cox and W. Josephson. File Synchronization with Vector Time Pairs. 2005.
- [12] Robert C. Daley and Peter G. Neumann. A general-purpose file system for secondary storage. pages 138–166, 2000.
- [13] David Dearman and Jeffery S. Pierce. “It’s on my other computer!”: computing with multiple devices. In Mary Czerwinski, Arnold M. Lund and Desney S. Tan (editors), *CHI*, pages 767–776, Florence, Italy, April 5-10 2008. ACM.

- [14] Ian Bicking. SQLObject. <http://www.sqlobject.org/>, 2008.
- [15] William Jones, Ammy Jiranida Phuwannurak, Rajdeep Gill and Harry Bruce. “Don’t take my folders away!”: organizing personal information to get things done. In Gerrit C. van der Veer and Carolyn Gale (editors), *CHI Extended Abstracts*, pages 1505–1508. ACM, 2005.
- [16] David R. Karger and William Jones. Data unification in personal information management. *Commun. ACM*, Volume 49, Number 1, pages 77–82, 2006.
- [17] DR Karger, W. Jones, O. Bergman, W. Pratt and M. Franklin. Towards a unification & integration of PIM support. *Jones, W., Bruce, H.: A Report on the NSF-Sponsored Workshop on Personal Information Management, Seattle, WA*, Volume 2005, pages 30–33, 2005.
- [18] Matt Mackall. Mercurial. <http://www.selenic.com/mercurial/wiki/>.
- [19] M. Mahalingam, C. Tang and Z. Xu. Towards a semantic, deep archival file system. *Distributed Computing Systems, 2003. FTDCS 2003. Proceedings. The Ninth IEEE Workshop on Future Trends of*, pages 115–121, 2003.
- [20] Brent Rector. Introducing Longhorn for Developers, Chapter 4: Storage. <http://msdn.microsoft.com/en-us/library/aa479870.aspx>.
- [21] Dennis M. Ritchie and Ken Thompson. The UNIX time-sharing system. *Commun. ACM*, Volume 17, Number 7, pages 365–375, 1974.
- [22] Craig A. N. Soules and Gregory R. Ganger. Connections: using context to enhance file search. *SIGOPS Oper. Syst. Rev.*, Volume 39, Number 5, pages 119–132, 2005.
- [23] The Coding Monkeys. SubEthaEdit. <http://www.codingmonkeys.de/subethaedit>, 2008.

Appendix A

Kernel Extension

The kernel extension developed for this project is known as **deveye**, and aims to coalesce various sources of raw data from within Mac OS X into a common interface accessible by any process running in userspace. This raw data represents the reporting of system calls as well as some particular interpretations of system calls. For instance, **deveye** is able to report a `close()` that occurs on any particular file, but only if that file has been modified during its time open (`SFE_MODIFIED`).

API

Installation

deveye is distributed as a single Mac OS X kernel extension, `deveye.kext`, along with a configuration file for the Mac OS X system manager, LaunchD. The kernel extension should be placed in `/System/Library/Extensions`, and the owner of the extension and its contents should be changed to `'root:wheel'`. Lastly, the configuration file, `com.moofco.deveye.plist` should be placed in `/Library/LaunchDaemons`. It should be loaded into the system by invoking `sudo launchctl load /Library/LaunchDaemons/com.moofco.deveye.plist`.

This will immediately start **deveye**, and reveal its public interface at `/dev/eyefs`. This API may now be used by user-space applications.

Usage

This section describes the interface for **deveye**. First, it is worth noting that `event.h` and `deveye.h`, the public header files for **deveye**, are available following. However, a short usage example is included here regardless.

Generally, an application must first `open()` the **deveye** device node at `/dev/eyefs`. A standard request structure should then be created, and allocated a static buffer as an output

target, before the entire request passed to `ioctl()`. A full example, `eyerender`, which simply displays all events generated within the kernel extension, is included on the following page.

com.moofco.deveye.plist

```
1 <?xml version='1.0' encoding='UTF-8'?>
2 <!DOCTYPE plist PUBLIC "-//Apple Computer//DTD PLIST 1.0//EN
3 http://www.apple.com/DTDs/PropertyList-1.0.dtd >
4 <plist version='1.0'>
5 <dict>
6 <key>Label</key>
7 <string>com.moofco.deveye</string>
8 <key>Debug</key>
9 <false />
10 <key>ProgramArguments</key>
11 <array>
12 <string>/System/Library/Extensions/deveye.kext/eyemgr</string>
13 </array>
14 <key>KeepAlive</key>
15 <true />
16 <key>RunAtLoad</key>
17 <true />
18 </dict>
19 </plist>
```

eyerender.c

```
1
2
3
4
5
6 #include <errno.h>
7 #include <fcntl.h>
8 #include <grp.h>
9 #include <libproc.h>
10 #include <pthread.h>
11 #include <pwd.h>
12 #include <stdint.h>
13 #include <stdio.h>
14 #include <stdlib.h>
15 #include <string.h>
16 #include <sys/ioccom.h>
17 #include <sys/ioctl.h>
18 #include <sys/param.h>
19 #include <sys/stat.h>
20 #include <sys/sysctl.h>
21 #include <sys/types.h>
22 #include <sys/uid.h>
23 #include <unistd.h>
24 #include "event.h"
25 #include "deveye.h"
26
27 int get_eyefs();
28
29 /**
30  * Core eyedm code.
31  */
32 int main() {
33
34     uint8_t output[SFE_MAXLEN];
35     int eyefs, status;
```

```
36 struct eyefs request;
37 struct sfe *outputp = ( struct sfe * ) output;
38 struct argsfe *arg;
39
40 memset( ( void * ) &request, 0, sizeof( struct eyefs ) );
41
42 // grab eyefs fd
43 if( ( eyefs = get_eyefs() ) < 0 ) {
44     fprintf( stderr, "couldn't open eyefs (%d) - is the kext loaded?\n", eyefs );
45 }
46
47 // poll for events from eyefs
48 for( ;; ) {
49
50     request.size = SFE_MAXLEN;
51     request.version = EYEFS_VERSION;
52     request.store.ptr = ( void * ) output;
53
54     if( status = ioctl( eyefs, EYEFS_GET, &request ) ) {
55         break;
56     }
57
58     printf( "event (last=%lld):\n", request.last );
59     printf( " -- uid=%d eid=%lld args=%d", outputp->uid, outputp->eid,
60           SFE_NUMARGS( outputp->event ) );
61     if( outputp->event & SFE_MOD ) {
62         printf( " MOD" );
63     }
64     if( outputp->event & SFE_FSEVENTS ) {
65         printf( " FSEVENTS" );
66     }
67     printf( "\n" );
68
69     if( arg = SFE_ARGO( outputp ) ) {
70         printf( " -- arg0 dev=%d uuid=%llu.%llu ino=%d\n",
71               arg->dev, *( ( uint64_t * ) &arg->uuid ),
72               *( ( uint64_t * ) ( ( char * ) &arg->uuid + 8 ) ),
73               arg->ino );
74         printf( " %s\n", arg->path );
75     }
76
77     if( arg = SFE_ARG1( outputp ) ) {
78         printf( " -- arg1 dev=%d uuid=%llu.%llu ino=%d\n",
79               arg->dev, *( ( uint64_t * ) &arg->uuid ),
80               *( ( uint64_t * ) ( ( char * ) &arg->uuid + 8 ) ),
81               arg->ino );
82         printf( " %s\n", arg->path );
83     }
84 }
85
86 // fall out
87 fprintf( stderr, "broke out: status=%d, errno=%d\n", status, errno );
88 return status;
89
90 }
91
92 /**
93  * Connect to the eyedm device node.
94  */
95 int get_eyefs() {
96
97
```

```

98 // is it really not more complicated than this?
99 return open( "/dev/" EYEFS_DEV, 0 );
100
101 }

```

event.h

```

1
2 /**
3  * event.h - provides the public and private (kernel) definitions for the sfe kext
4  */
5 #ifndef __EVENT_H
6 #define __EVENT_H
7
8 #include <sys/ioccom.h>
9 #include <sys/param.h>
10 #include <sys/vnode.h>
11 #include <sys/kernel.h>
12 #include <stdint.h>
13
14 /**
15  * UUID typedef - typically included out of kernel, but not available in kernel.
16  */
17 #ifdef KERNEL
18 typedef unsigned char uuid_t[16];
19 #endif
20
21 /**
22  * Header for a single dynamic-length event (24 bytes).
23  */
24 struct sfe {
25     uint64_t eid;
26     uid_t uid;
27     gid_t gid;
28     pid_t pid;
29     uint32_t event;
30 };
31
32 /**
33  * Arguments for the event (24 bytes). The actual size of this struct will be
34  * inclusive of the path length, which should be 8-byte aligned.
35  */
36 struct argsfe {
37     uuid_t uuid;
38     dev_t dev;
39     ino_t ino;
40     uid_t uid;
41     gid_t gid;
42     int32_t mode, length;
43     char path[0];
44 };
45
46 /**
47  * Internal helper methods for the usage of events.
48  */
49 #ifdef KERNEL
50 int eventsetup();
51 int eventshutdown();
52 int uuid_submit( struct argsfe * );
53 int uuid_fetch( struct argsfe * );

```

```

54 int queue_push( struct sfe * );
55 int queue_push_user( user_addr_t, int );
56 int event_vnode( kauth_cred_t, kauth_action_t, vfs_context_t, vnode_t, vnode_t );
57 int event_fileop( kauth_cred_t, kauth_action_t, vfs_context_t, vnode_t );
58 int event_obtain( uid_t, uint64_t, uint8_t, struct sfe * );
59 int event_obtain_user( uid_t, uint64_t *, uint8_t, user_addr_t, int * );
60 #endif
61
62 /**
63  * The maximum length of any single event.
64  */
65 #define SFE_MAXLEN ( sizeof( struct sfe ) + 2 * ( sizeof( struct argsfe ) + MAXPATHLEN ) )
66
67 /**
68  * Definitions to return the total size of a whole event along with the total
69  * size of a single argument.
70  */
71 #define SFE_ARGSIZE(x) ( x ? (sizeof(struct argsfe)+((struct argsfe*)x)->length) : 0 )
72 #define SFE_ARGO(v) ( \
73     ( (v)->event & SFE_ARG ) == SFE_ARG ? \
74     ((struct argsfe*)((void*)(v))+sizeof(struct sfe)) : 0 \
75 )
76 #define SFE_ARG1(v) ( \
77     ( (v)->event & SFE_ARGS ) == SFE_ARGS ? \
78     ((struct argsfe*)((void*)SFE_ARGO(v))+sizeof(struct argsfe)+SFE_ARGO(v)->length) : 0 \
79 )
80 #define SFE_SIZE(v) ( \
81     ( ( (v)->event & SFE_ARG ) == SFE_ARG ? SFE_ARGSIZE(SFE_ARGO(v)) : 0 ) + \
82     ( ( (v)->event & SFE_ARGS ) == SFE_ARGS ? SFE_ARGSIZE(SFE_ARG1(v)) : 0 ) + \
83     ( sizeof( struct sfe ) ) \
84 )
85
86 /**
87  * Define high-byte properties (part of events).
88  */
89 #define SFE_MOD        0x0100
90 #define SFE_ARG        0x0200
91 #define SFE_ARGS       0x0600
92 #define SFE_MERGED     0x1000
93 #define SFE_FSEVENTS   0x2000
94 #define SFE_PRIVATE    0x4000
95
96 /**
97  * Generates the number of arguments from an event field.
98  */
99 #define SFE_NUMARGS(v) ((v & SFE_ARGS) == SFE_ARGS ? 2 : ((v & SFE_ARG) ? 1 : 0))
100
101 /**
102  * Event type definitions. The event field may be masked with 0xff to find
103  * specific event types.
104  */
105 #define SFE_OPEN        ( 0x01 | SFE_ARG ) // open notification
106 #define SFE_READ        ( 0x02 | SFE_ARG ) // authorisation to read
107 #define SFE_EXEC        ( 0x03 | SFE_ARG ) // authorisation to execute/browse
108 #define SFE_WRITE       ( 0x04 | SFE_ARG ) // authorisation to write (i.e. no changes have been made yet)
109 #define SFE_CREATE      ( 0x05 | SFE_MOD | SFE_ARG )
110 #define SFE_DELETE      ( 0x06 | SFE_MOD | SFE_ARG )
111 #define SFE_STAT        ( 0x07 | SFE_MOD | SFE_ARG ) // file stats changed
112 #define SFE_RENAME      ( 0x08 | SFE_MOD | SFE_ARGS ) // if the second argument has a uid/gid, it has been replaced
113 #define SFE_MODIFIED    ( 0x09 | SFE_MOD | SFE_ARG ) // dirty close
114 #define SFE_EXCHANGE    ( 0x0a | SFE_MOD | SFE_ARGS )
115 #define SFE_FINDER      ( 0x0b | SFE_MOD | SFE_ARG ) // finder info changed

```

```

116 #define SFE_MKDIR      ( 0x0c | SFE_MOD | SFE_ARG )
117 #define SFE_CHOWN     ( 0x0d | SFE_MOD | SFE_ARG )
118 #define SFE_MOUNT     ( 0x0e | SFE_PRIVATE | SFE_ARG )
119
120 #endif/*_EVENT_H*/

```

deveye.h

```

1
2 /**
3  * deveye.h - provides simple definitions for user-space to communicate with us
4  */
5 #ifndef __DEVEYE_H
6
7 #include <sys/ioccom.h>
8 #include <sys/types.h>
9 #include <stdint.h>
10
11 /**
12  * Global version, device name created in /dev, along with xattr definitions.
13  */
14 #define EYEFS_VERSION    2
15 #define EYEFS_DEV       "eyefs"
16 #define EYEFS_TLD       "com.moofco.deveye"
17 #define EYEFS_UUID      "uuid"

```

```

18
19 /**
20  * Request class for access to our standard device (EYEFS_DEV). Includes an
21  * internal union for 32-64 transition between the real world and kernel land.
22  */
23 struct eyefs {
24     int16_t version, count;
25     int32_t size;
26     uint64_t last;
27     union {
28         void *ptr;
29         user_addr_t addr;
30         int32_t old_addr;
31     } store;
32 };
33
34 /**
35  * Device ioctl() request actions.
36  * - GET: waits and returns events
37  * - POLL: same as GET, but will return immediately if there are no pending events
38  * - PUSH: used internally by 'eyemgr' to submit fsevents
39  */
40 #define EYEFS_GET        _IOWR( 's', 1, struct eyefs )
41 #define EYEFS_POLL      _IOWR( 's', 2, struct eyefs )
42 #define EYEFS_PUSH      _IOW( 's', 3, struct eyefs )
43
44 #endif/*_DEVEYE_H*/

```

Internals

We have also included the internal source code for `deveye`. The vast majority of the init/configuration code is within `deveye.c`, with much of the conceptual storage and list code offloaded into `event.c`.

deveye.c

```
1
2 /**
3  * deveye.c - provides the core outwards controls/functions for the next.
4  */
5
6 #include "deveye.h"
7 #include "event.h"
8
9 #include <libkern/OSAtomic.h>
10 #include <libkern/libkern.h>
11 #include <miscfs/devfs/devfs.h>
12 #include <sys/conf.h>
13 #include <sys/kauth.h>
14 #include <sys/kernel.h>
15 #include <sys/system.h>
16 #include <sys/unistd.h>
17 #include <sys/vnode.h>
18 #include <sys/errno.h>
19
20 static kauth_listener_t gListener = NULL;
21 static SInt32 gActivationCount = 0;
22 static SInt32 gDevInstalled = -1;
23 static void *gDevNode = NULL;
24
25 /**
26  * Standard device switch open/close implementation.
27  */
28 static int dev_opcl(
29     __unused dev_t dev,
30     __unused int flag,
31     __unused int mode,
32     __unused struct proc *p ) {
33     return 0;
34 }
35
36 /**
37  * Handles all user-space calls to ioctl, providing events either as a
38  * 'get' or 'poll' operation.
39  */
40 static int dev_ioctl(
41     __unused dev_t dev,
42     u_long cmd,
43     caddr_t data,
44     __unused int flag,
45     struct proc *p ) {
46
47     int status = EIO;
48     kauth_cred_t cred;
49     uid_t uid;
50     struct eyeifs *request = ( struct eyeifs * ) data;
51     user_addr_t target;
52
53     // discard invalid requests
54     if( request->version != EYEFS_VERSION ) {
55         return ENOPROTOPT; // protocol not available
56     }
57
58     // increment usage counter, grab cred/uid for validation
59     OSIncrementAtomic( &gActivationCount );
60     cred = kauth_cred_proc_ref( p );
61
62     uid = kauth_cred_getuid( cred );
63
64     // all a bit strange, but seems to work . . .
65     if( !proc_is64bit( p ) ) {
66         target = request->store.addr;
67     }
68     else {
69         // target = CAST_USER_ADDR_T( request->store.old_addr );
70         target = request->store.addr;
71     }
72
73     switch( cmd ) {
74     case EYEFS_GET:
75     case EYEFS_POLL:
76         // submit request internally
77         status = event_obtain_user( uid, &( request->last ), cmd == EYEFS_POLL, target, &( request->size ) );
78         if( status ) {
79             request->count = 0;
80         }
81         else {
82             request->count = 1; // for now, we only have a single result
83         }
84         break;
85     case EYEFS_PUSH:
86         // this may only be invoked in push
87         if( uid ) {
88             status = EACCES;
89             goto fallout;
90         }
91
92         // FIXME: not yet supported
93         if( request->count != 1 ) {
94             status = E2BIG;
95             goto fallout;
96         }
97
98         // submit request internally
99         status = queue_push_user( target, request->size );
100         break;
101     }
102
103 fallout:
104     // release resources, return determined status
105     kauth_cred_unref( &cred );
106     OSDecrementAtomic( &gActivationCount );
107     return status;
108 }
109
110 /**
111  * Device switch structure, points only to default open/close constructs, and our major ioctl implementation.
112  */
113 static struct cdevsw gDevSwitch = {
114     dev_opcl, dev_opcl,
115     eno_rdwrt, eno_rdwrt, dev_ioctl, eno_stop, eno_reset,
116     NULL, eno_select, eno_mmap, eno_strat, eno_getc, eno_putc, 0
117 };
118
119 /**
120  * The core kauth handler for vnode operations.
121  */
122
```

```

123  */
124  static int kauth_fileop( kauth_cred_t cred, __unused void *idata,
125      kauth_action_t action,
126      uintptr_t arg0, uintptr_t arg1, uintptr_t arg2, uintptr_t arg3 ) {
127
128      vnode_t vp = ( vnode_t ) arg0;
129      vfs_context_t ctx;
130      int uid;
131
132      OSIncrementAtomic( &gActivationCount );
133      ctx = vfs_context_create( ( vfs_context_t ) 0 );
134
135      if( action == KAUTH_FILEOP_OPEN && ( uid = kauth_cred_getuid( cred ) ) ) {
136          if( event_fileop( cred, action, ctx, vp ) ) {
137              // do we care?
138          }
139      } else if( action == KAUTH_FILEOP_EXEC && ( uid = kauth_cred_getuid( cred ) ) ) {
140          if( event_fileop( cred, action, ctx, vp ) ) {
141              // do we care?
142          }
143      }
144
145      vfs_context_rele( ctx );
146      OSDecrementAtomic( &gActivationCount );
147
148      return KAUTH_RESULT_DEFER;
149
150  }
151  /**
152   * Core kernel extension init function.
153   */
154  kern_return_t deveye_start( kmod_info_t *ki, void *d ) {
155
156      // setup internal event storage
157      if( eventsetup() ) {
158          printf( "deveye_start: could not eventsetup()" );
159          return KERN_FAILURE;
160      }
161
162      // create the kauth listener.
163      gListener = kauth_listen_scope( KAUTH_SCOPE_FILEOP, kauth_fileop, NULL );
164      if( gListener == NULL ) {
165          printf( "deveye_start: could not create gListener for vnodes.\n" );
166          return KERN_FAILURE;
167      }
168
169      // create a new device node
170      gDevInstalled = cdevsw_add( -1, &gDevSwitch );
171      if( gDevInstalled < 0 ) {
172          printf( "deveye_start: could not insert cdevsw\n" );
173          return KERN_FAILURE;
174      }
175
176      // install the device node in '/dev'
177      gDevNode = devfs_make_node( makedev( gDevInstalled, 0 ),
178          DEVFS_CHAR, UID_ROOT, GID_WHEEL, 0644, EYEFS_DEV, 0 );
179      if( NULL == gDevNode ) {
180          printf( "deveye_start: could not make device node\n" );
181          return KERN_FAILURE;
182      }
183
184      // success!

```

```

185      printf( "deveye_start: ok\n" );
186      return KERN_SUCCESS;
187  }
188
189  /**
190   * Core shutdown function.
191   */
192  kern_return_t deveye_stop( kmod_info_t *ki, void *d ) {
193
194      // revoke event setup (has to happen first - kicks the habit of listeners)
195      eventshutdown();
196
197      // wait for our activation count to drop
198      do {
199          struct timespec oneSecond;
200          oneSecond.tv_sec = 1;
201          oneSecond.tv_nsec = 0;
202          msleep( &gActivationCount, NULL, PUSER, "com.moofco.kext.deveye.stop", &oneSecond );
203      } while ( gActivationCount > 0 );
204
205      // release the kauth listener
206      if( gListener != NULL ) {
207          kauth_unlisten_scope( gListener );
208          gListener = NULL;
209      }
210
211      // release the device node
212      if( gDevInstalled > -1 ) {
213          cdevsw_remove( gDevInstalled, &gDevSwitch );
214          gDevInstalled = -1;
215      }
216
217      // uninstall the device node from '/dev'
218      if( gDevNode != NULL ) {
219          devfs_remove( gDevNode );
220          gDevNode = NULL;
221      }
222
223      // success
224      printf( "deveye_stop: ok\n" );
225      return KERN_SUCCESS;
226  }
227
228  }

```

event.c

```

1
2  /**
3   * event.c - provides tighter-knit functionality for managing events
4   */
5
6  #include <libkern/OSBase.h>
7  #include <libkern/OSAtomic.h>
8  #include <libkern/libkern.h>
9  #include <miscfs/devfs/devfs.h>
10 #include <sys/conf.h>
11 #include <sys/kauth.h>
12 #include <sys/kernel.h>
13 #include <sys/system.h>

```

```

14 #include <sys/unistd.h>
15 #include <sys/vnode.h>
16 #include <sys/errno.h>
17 #include "fastlock.h"
18 #include "event.h"
19
20 #ifdef DEBUG
21 #define dprintf(format, args...) printf(format, ## args)
22 #else
23 #define dprintf(format, args...) 0
24 #endif
25
26 #define QUEUEBITS 13
27 #define QUEUELEN (1<<QUEUEBITS)
28 #define QUEUEMASK (QUEUELEN-1)
29
30 static SInt64 gUniqueEvent = 0;
31 static SInt32 gQueuePos = 0;
32 static SInt32 gQueueActive = 0;
33 static struct {
34     fastlock_t lck;
35     union {
36         struct sfe sfe;
37         uint8_t __internal[SFE_MAXLEN];
38     } ev;
39 } gQueue[QUEUELEN];
40
41 #define MAPBITS 13
42 #define MAPLEN (1<<MAPBITS)
43
44 static struct {
45     uuid_t uuid;
46     dev_t dev;
47 } gDevUUID[MAPLEN];
48
49 /**
50  * Configure the event queue structure.
51  */
52 int eventsetup() {
53
54     memset( gQueue, 0, sizeof( gQueue ) );
55     memset( gDevUUID, 0, sizeof( gDevUUID ) );
56
57     gUniqueEvent = 0;
58     gQueuePos = 0;
59     gQueueActive = 1;
60
61     printf( "deveye: eventsetup using %d bytes for queue,"
62           " %d bytes for dev/uuid (now open for business)\n",
63           sizeof( gQueue ), sizeof( gDevUUID ) );
64
65     return 0;
66 }
67
68 /**
69  * Notify our listeners that the queue structure is shutting down.
70  */
71 int eventshutdown() {
72
73     gQueueActive = 0;
74     wakeup( ( caddr_t ) &gQueue );
75

```

43

```

76
77     printf( "deveye: eventshutdown ok\n" );
78
79     return 0;
80
81 }
82
83 /**
84  * Hash a device node.
85  */
86 int uuid_devhash( dev_t dev ) {
87
88     int a = ( int ) dev;
89
90     a = a ^ ( a >> 4 );
91     a = ( a ^ 0xdeadbeef ) + ( a << 5 );
92     a = a ^ ( a >> 11 );
93     a = a % MAPLEN;
94
95     return a;
96 }
97
98 /**
99  * Device vs. UUID push.
100  */
101 int uuid_submit( struct argsfe *arg ) {
102
103     int i, target = uuid_devhash( arg->dev );
104
105 #ifdef DEBUG
106     // ensure that this method has been called properly
107     if( !arg ) {
108         printf( "deveye: uuid_submit without arg\n" );
109         return -2;
110     }
111     dprintf( "deveye: trying to place dev %d into bucket %d\n", arg->dev, target );
112 #endif
113
114     // find a bucket that this uuid/dev can live in
115     for( i = target + 1; ( i % MAPLEN ) != target; ++i ) {
116         if( gDevUUID[i % MAPLEN].dev == 0 || gDevUUID[i % MAPLEN].dev == arg->dev ) {
117             break;
118         }
119     }
120     i %= MAPLEN;
121
122     // check rollover conditions
123     if( i == target ) {
124         printf( "deveye: couldn't fit device uuid for dev %d\n", arg->dev );
125         return ENOMEM;
126     }
127
128     // place the uuid where it should be
129     gDevUUID[i].dev = arg->dev;
130     memcpy( &( gDevUUID[i].uuid ), &( arg->uuid ), 16 );
131
132     // success!
133     printf( "deveye: dev %d[%d] given uuid %llu.%llu\n", arg->dev, i, *( ( uint64_t * ) &arg->uuid ), *( ( uint64_t * ) ( ( ( char *
134     return 0;
135
136
137

```



```

138 }
139
140 /**
141  * Copy out a UUID if its found.
142  */
143 int uuid_fetch( struct argsfe *arg ) {
144
145     int i, target = uuid_devhash( arg->dev ), status = 0;
146
147 #ifndef DEBUG
148     // ensure that this method has been called properly
149     if( !arg ) {
150         printf( "deveye: uuid_fetch without arg\n" );
151         return -2;
152     }
153     dprintf( "deveye: trying to look for dev %d in bucket %d\n", arg->dev, target );
154 #endif
155
156     // find the bucket containing this dev
157     for( i = target + 1; ( i % MAPLEN ) != target; ++i ) {
158         if( gDevUUID[i % MAPLEN].dev == arg->dev ) {
159             break;
160         }
161         if( gDevUUID[i % MAPLEN].dev == 0 ) {
162             status = EIO;
163             goto fallout;
164         }
165     }
166     i %= MAPLEN;
167
168     // check rollover conditions (unlikely, but eh)
169     if( i == target ) {
170         status = ENOMEM;
171         goto fallout;
172     }
173
174     // copy the uuid into the output struct
175     memcpy( &( arg->uuid ), &( gDevUUID[i].uuid ), 16 );
176
177     // success!
178     dprintf( "deveye: found dev %d in bucket %d\n", arg->dev, i );
179     return 0;
180
181 fallout:
182     // fallout condition, reset the uuid to blank
183     memset( &( arg->uuid ), 0, 16 );
184     return status;
185 }
186
187 /**
188  * Internal (single) event push, from kernel memory.
189  */
190 int queue_push( struct sfe *ev ) {
191
192     int32_t pos;
193
194     // handle private events
195     if( ev->event & SFE_PRIVATE ) {
196         int status = ENOTSUP;
197
198         switch( ev->event ) {

```

```

200         case SFE_MOUNT:
201             status = uuid_submit( SFE_ARGO( ev ) );
202             break;
203     }
204
205     return status;
206 }
207
208 // update event, eid and uuids for args
209 ev->eid = OSAddAtomic64( 1, &gUniqueEvent );
210 if( ( ev->event & SFE_ARG ) == SFE_ARG ) {
211     uuid_fetch( SFE_ARGO( ev ) );
212     if( ( ev->event & SFE_ARGS ) == SFE_ARGS ) {
213         uuid_fetch( SFE_ARG1( ev ) );
214     }
215 }
216
217 // grab next position, copy, free lock
218 for( ;; ) {
219     pos = OSIncrementAtomic( ( SInt32 * ) &gQueuePos ) & QUEUEMASK;
220     if( FastSIntLockEx( &gQueue[pos].lck ) ) {
221         break;
222     }
223     printf( "deveye: could not lock at %d to insert, trying next.\n", pos );
224 }
225 memcpy( ( void * ) &gQueue[pos].ev, ( void * ) ev, SFE_SIZE( ev ) );
226 FastSIntUnlockEx( &gQueue[pos].lck );
227
228 // wakeup listeners
229 wakeup( ( caddr_t ) &gQueue );
230
231 // success
232 return 0;
233 }
234
235 /**
236  * Inserts a single SFE into the kernel.
237  */
238 int queue_push_user( user_addr_t ev, int length ) {
239
240     uint8_t buffer[SFE_MAXLEN];
241
242     // confirm we've been given enough data
243     if( length < sizeof( struct sfe ) ) {
244         printf( "deveye: copyin not valid - len=%d, len(sfe)=%d\n", length, sizeof( struct sfe ) );
245         return EINVAL;
246     }
247
248     // copy in from user-space
249     if( copyin( ev, ( void * ) buffer, length ) ) {
250         printf( "deveye: copyin failed - len=%d\n", length );
251         return EFAULT;
252     }
253
254     // push the event
255     return queue_push( ( struct sfe * ) buffer );
256 }
257
258 /**
259  * Does the internal work regarding the creation of a formal event from KAuth-vnode data.
260  */

```

```

262  */
263  int event_fileop( kauth_cred_t cred, kauth_action_t action, vfs_context_t ctx, vnode_t vp ) {
264
265      uint8_t buffer[SFE_MAXLEN];
266      struct sfe *ev = ( struct sfe * ) buffer;
267      struct argsfe *arg = ( struct argsfe * ) ( buffer + sizeof( struct sfe ) );
268      struct vnode_attr va;
269
270      // grab standard vars
271      ev->uid = kauth_cred_getuid( cred );
272      ev->gid = kauth_cred_getgid( cred );
273      ev->pid = proc_selfpid();
274
275      // set up vnode request
276      VATTR_INIT( &va );
277      VATTR_WANTED( &va, va_fsid );
278      VATTR_WANTED( &va, va_fileid );
279      VATTR_WANTED( &va, va_mode );
280      VATTR_WANTED( &va, va_uid );
281      VATTR_WANTED( &va, va_gid );
282
283      // throw out if this vnode can't be found
284      if( vnode_getattr( vp, &va, ctx ) ) {
285          return -1;
286      }
287
288      // grab details from vnode
289      memset( arg->uuid, 0, 16 );
290      arg->dev = ( dev_t ) va.va_fsid;
291      arg->ino = ( ino_t ) va.va_fileid;
292      arg->mode = ( int32_t ) VTOIF( vnode_vtype( vp ) ) | va.va_mode;
293      arg->uid = ( uid_t ) va.va_uid;
294      arg->gid = ( gid_t ) va.va_gid;
295
296      // grab path from vnode, and store it aligned to a 8-byte word
297      arg->length = MAXPATHLEN;
298      if( vn_getpath( vp, arg->path, &( arg->length ) ) ) {
299          arg->path[0] = '\0';
300      }
301      arg->length = ( ( strlen( arg->path ) >> 3 ) + 1 ) << 3;
302
303      // submit fileop open
304      if( action == KAUTH_FILEOP_OPEN ) {
305          ev->event = SFE_OPEN;
306          queue_push( ev );
307      }
308  /*
309      // submit vnode auth read
310      if( action & KAUTH_VNODE_READ_DATA ) {
311          ev->event = SFE_READ;
312          queue_push( ev );
313      }
314  */
315
316      // submit vnode auth execute
317      if( action == KAUTH_FILEOP_EXEC ) {
318          ev->event = SFE_EXEC;
319          queue_push( ev );
320      }
321
322  /*
323      // submit vnode auth write
324      if( action & KAUTH_VNODE_WRITE_DATA ) {
325          ev->event = SFE_WRITE;
326          queue_push( ev );
327      }
328  */
329      // done and done
330      return 0;
331  }
332
333  /**
334   * Finds the event with id most closest (greater) to 'last' (or any, if zero).
335   */
336  int event_obtain( uid_t uid, uint64_t last, uint8_t poll, struct sfe *ev ) {
337
338      uint64_t event_start;
339      int32_t pos, closest, count;
340
341  retry:
342      count = 0;
343      closest = -1;
344      pos = OSAddAtomic( 0, &QueuePos );
345      event_start = OSAddAtomic64( 0, &UniqueEvent );
346
347      // check shutdown condition
348      if( !gQueueActive ) {
349          return EBUSY;
350      }
351
352      // work backwards
353      for( ; ; ) {
354          pos = ( pos - 1 ) & QUEUEMASK;
355
356          // check raw wrap
357          if( ++count == QUEUELEN ) {
358              printf( "deveye: had a simple wrap\n" );
359              pos = -1;
360              break;
361          }
362
363          // try to get a shared lock on this object - if we can't, we've wrapped (or, uh, just started)
364          if( FastSintLock( &Queue[pos].lck ) < 0 ) {
365              dprintf( "deveye: pos %d was being written to\n", pos );
366              continue;
367          }
368
369          // check other wrap conditions (blank entry, id wrap)
370          if( !gQueue[pos].ev.sfe.eid || gQueue[pos].ev.sfe.eid > event_start ) {
371              FastSintUnlock( &Queue[pos].lck );
372              printf( "deveye: pos %d lock ok, but was blank or a wrap\n", pos );
373              pos = -1;
374              break;
375          }
376      }
377
378      // check uid condition
379      if( !uid || gQueue[pos].ev.sfe.uid == uid ) {
380
381          // if we don't care what event we find
382          if( !last ) {
383              dprintf( "deveye: got pos %d because we didn't care\n", pos );
384              break;
385          }

```

```

386
387 // if we're not quite there yet, store closest
388 if( gQueue[pos].ev.sfe.eid > last ) {
389     if( closest != -1 ) {
390         FastSIntUnlock( &gQueue[closest].lck );
391     }
392     closest = pos;
393     continue;
394 }
395
396 // if we've hit our limit - note we don't actually want /this/ event
397 if( gQueue[pos].ev.sfe.eid <= last ) {
398     FastSIntUnlock( &gQueue[pos].lck );
399     dprintf( "deveye: hit our eid limit on pos %d\n", pos );
400     pos = -1;
401     break;
402 }
403
404 }
405
406 // unlock before loop
407 FastSIntUnlock( &gQueue[pos].lck );
408
409 }
410
411 // if we had focus on a closest ok element, use that
412 if( closest != -1 ) {
413     pos = closest;
414     dprintf( "deveye: got pos %d - was closest to eid limit\n", pos );
415 }
416
417 // check outcome of search
418 if( pos == -1 ) {
419
420     // ignore if we're polling
421     if( poll ) {
422         return EAGAIN;
423     }
424
425     // check if we should wrap without waiting for a signal
426     if( event_start != OSAddAtomic64( 0, &gUniqueEvent ) ) {
427         dprintf( "deveye: we have new events,"
428             " we don't need to bother sleeping\n" );
429         goto retry;
430     }
431
432     // sleep for a new event
433     dprintf( "deveye: sleeping for relevant event\n" );
434     if( msleep( ( caddr_t ) &gQueue, NULL, PCATCH, "event_obtain", NULL ) ) {
435         return EINTR;
436     }
437     goto retry;
438
439 }
440
441 // return this event
442 dprintf( "deveye: returning event at %d\n", pos );
443 memcpy( ( void * ) ev,
444     ( void * ) &( gQueue[pos].ev.sfe ),
445     SFE_SIZE( &( gQueue[pos].ev.sfe ) ) );
446 FastSIntUnlock( &gQueue[pos].lck );
447

```

```

448 // success
449 return 0;
450
451 }
452
453 /**
454  * Quick hack to insert a single SFE into the kernel.
455  */
456 int event_obtain_user(
457     uid_t uid, uint64_t *last, uint8_t poll, user_addr_t ev, int *length ) {
458
459     uint8_t buffer[SFE_MAXLEN];
460     int status, size;
461
462     // confirm we've been given enough space to work with
463     if( *length < SFE_MAXLEN ) {
464         printf( "deveye: copyout might fail - len=%d, len(sfe)=%d\n",
465             *length, sizeof( struct sfe ) );
466         return EINVAL;
467     }
468
469     // find a relevant event
470     if( status = event_obtain( uid, *last, poll, ( struct sfe * ) buffer ) ) {
471         return status;
472     }
473
474     // copyout a single event
475     size = SFE_SIZE( ( struct sfe * ) buffer );
476     if( copyout( ( void * ) buffer, ev, size ) ) {
477         printf( "deveye: copyout failed - len=%d\n", size );
478         return EFAULT;
479     }
480
481     // update last event pointer
482     *last = ( ( struct sfe * ) buffer )->eid;
483     return 0;
484
485 }

```

eyemgr.c

This tool is compiled and run as root, supporting the kernel extension directly with information sourced from the private FSEvents [2] API.

```

1
2 /**
3  * eyebmgr.c - handles module load, and fsevents hook
4  * in the not-too distant future, it would be great to have this all inside the kernel, but for now - here it is
5  */
6
7 #include <CoreServices/CoreServices.h>
8 #include <CoreFoundation/CoreFoundation.h>
9 #include <errno.h>
10 #include <fcntl.h>
11 #include <grp.h>
12 #include <libproc.h>
13 #include <pthread.h>
14 #include <pwd.h>
15 #include <stdint.h>
16 #include <stdio.h>

```

```

17 #include <stdlib.h>
18 #include <string.h>
19 #include <sys/ioccom.h>
20 #include <sys/ioctl.h>
21 #include <sys/param.h>
22 #include <sys/mount.h>
23 #include <sys/stat.h>
24 #include <sys/sysctl.h>
25 #include <sys/types.h>
26 #include <sys/unistd.h>
27 #include <sys/uio.h>
28 #include <sys/xattr.h>
29 #include <unistd.h>
30 #include <pthread.h>
31
32 #include "fsevents.h"
33 #include "event.h"
34 #include "deveye.h"
35
36 #define BUFLen 1024*16
37
38 int get_fsevents();
39 int get_eyefs();
40 int parse_fsevent( void *, void * );
41 uint32_t sfe_type( int32_t );
42 void register_volume_path( const char * );
43 void fsevents_callback(
44     ConstFSEventStreamRef,
45     void *, size_t, void *,
46     const FSEventStreamEventFlags eventFlags[],
47     const FSEventStreamEventId eventIds[] );
48 void *mount_mgr( void * );
49
50 static int eyefs;
51
52 /**
53  * Core eyemgr code.
54  */
55 int main() {
56
57     uint8_t buffer[BUFLen], output[SFE_MAXLEN];
58     int fsevents, length, parsed, status = EIO;
59     struct eyefs request;
60     struct sfe *outputp = ( struct sfe * ) output;
61 #ifdef DEBUG
62     struct argsfe *temparg = ( struct argsfe * ) ( output + sizeof( struct sfe ) );
63     struct argsfe *temparg2;
64 #endif
65     pthread_t th;
66
67     sranddev();
68
69     // set up eyefs request struct (won't change)
70     request.count = 1;
71     request.version = EYEFS_VERSION;
72     request.store.ptr = ( void * ) output;
73
74     // check that we are root
75     if( getuid() != 0 ) {
76         fprintf( stderr, "eyemgr; should be invoked as root"
77             " (quit fucking around :D)\n" );
78     }

```

```

79
80 // grab fsevents fd
81 if( ( fsevents = get_fsevents() ) < 0 ) {
82     fprintf( stderr, "eyemgr; couldn't open fsevents (%d)"
83         " - is this tiger or higher?\n", fsevents );
84 }
85
86 // load the deveye kext
87 #ifdef DEBUG
88     fprintf( stderr, "eyemgr; DEBUG - not attempting to load deveye kext\n" );
89 #else
90     if( system( "/sbin/kextload /System/Library/Extensions/deveye.kext" ) ) {
91         fprintf( stderr, "eyemgr; couldn't load/find deveye kext\n" );
92     }
93 #endif
94
95 // grab eyefs fd
96 if( ( eyefs = get_eyefs() ) < 0 ) {
97     fprintf( stderr, "eyemgr; couldn't open eyefs (%d)"
98         " - was the kext loaded ok?\n", eyefs );
99 }
100
101 // create mount/umount thread
102 pthread_create( &th, NULL, mount_mgr, NULL );
103
104 // core loop, read from fsevents
105 while( 0 < ( length = read( fsevents, buffer, sizeof( buffer ) ) ) ) {
106
107     // we might be given multiple events; parse them all
108     for( parsed = 0; parsed != length; ) {
109
110         // clear output struct (probably not required)
111         memset( ( void * ) output, 0, SFE_MAXLEN );
112
113         // set up request
114         parsed += parse_fsevent( buffer + parsed, output );
115         request.size = SFE_SIZE( outputp );
116
117         // FIXME: remove debug output
118 #ifdef DEBUG
119         printf( "got event (size=%d): ev=%d, arg0=(%d) %s",
120             request.size, outputp->event, temparg->length, temparg->path );
121         if( ( outputp->event & SFE_ARGS ) == SFE_ARGS ) {
122             temparg2 = ( struct argsfe * ) ( ( void * ) temparg )
123                 + sizeof( struct argsfe ) + temparg->length;
124             printf( ", arg1=(%d) %s", temparg2->length, temparg2->path );
125         }
126         printf( "\n" );
127 #endif
128
129         // submit to eyefs
130         if( ioctl( eyefs, EYEFS_PUSH, &request ) ) {
131             status = errno;
132             goto fallout;
133         }
134     }
135 }
136
137 }
138
139 fallout:
140 // fsevents is shutting down; machine is probably going off

```

```

141     if( length == 0 ) {
142         status = 0;
143     }
144     fprintf( stderr, "falling out: length=%d, status=%d\n", length, status );
145
146     // fall out
147     return status;
148
149 }
150
151 /**
152  * Connect to the fsevents device.
153  */
154 int get_fsevents() {
155
156     int fd, fsfd;
157     fsevent_clone_args retrieve_ioctl;
158     signed char event_list[FSE_MAX_EVENTS];
159
160     memset( event_list, 0, sizeof( event_list ) );
161
162     // set up event list
163     event_list[FSE_CREATE_FILE]      = FSE_REPORT;
164     event_list[FSE_DELETE]           = FSE_REPORT;
165     event_list[FSE_STAT_CHANGED]     = FSE_REPORT;
166     event_list[FSE_RENAME]           = FSE_REPORT;
167     event_list[FSE_CONTENT_MODIFIED] = FSE_REPORT;
168     event_list[FSE_EXCHANGE]         = FSE_REPORT;
169     event_list[FSE_FINDER_INFO_CHANGED] = FSE_REPORT;
170     event_list[FSE_CREATE_DIR]       = FSE_REPORT;
171     event_list[FSE_CHOWN]            = FSE_REPORT;
172
173     // set up ioctl request
174     retrieve_ioctl.event_list = event_list;
175     retrieve_ioctl.num_events = sizeof( event_list );
176     retrieve_ioctl.event_queue_depth = -1;
177     retrieve_ioctl.fd = &fd;
178
179     // open fsevents device, get inner fd
180     if( 0 > ( fsfd = open( "/dev/fsevents", 0 ) ) ) {
181         return -2;
182     }
183     if( 0 > ioctl( fsfd, FSEVENTS_CLONE, &retrieve_ioctl ) ) {
184         return -3;
185     }
186     close( fsfd );
187
188     // success
189     return fd;
190
191 }
192
193 /**
194  * Connect to the eyemgr device node.
195  */
196 int get_eyefs() {
197
198     // is it really not more complicated than this?
199     return open( "/dev/" EYEFS_DEV, 0 );
200
201 }
202

```

```

203 /**
204  * Creates a single sfe (in stream format) from the given fsevents input buffer.
205  */
206 int parse_fsevent( void *buffer, void *output ) {
207
208     struct sfe *ev = ( struct sfe * ) output;
209     struct argsfe *arg = NULL;
210     struct {
211         int32_t type;
212         pid_t pid;
213     } __attribute__((__packed__)) *header = ( void * ) buffer;
214     struct {
215         uint16_t type;
216         uint16_t length;
217         union {
218             char str[0];
219             dev_t dev;
220             ino_t inode;
221             int32_t mode;
222             uid_t uid;
223             gid_t gid;
224             int32_t int32;
225             int64_t int64;
226         } v;
227     } __attribute__((__packed__)) *fsarg;
228     struct proc_bsdinfo pinfo;
229     int file;
230
231     ev->event = sfe_type( header->type ) | SFE_FSEVENTS;
232     // identify this event as from fsevents
233     ev->pid = header->pid;
234     ev->uid = -1;
235     ev->gid = -1;
236     file = -1;
237
238     // attempt to grab the uid/gid of this event via system call
239     if( 0 < proc_pidinfo( header->pid, PROC_PIDTBSDINFO, 0, ( void * ) &pinfo, sizeof( struct proc_bsdinfo ) ) ) {
240         ev->uid = pinfo.pbi_uid;
241         ev->gid = pinfo.pbi_gid;
242     }
243
244     // parse every argument to this fsevent
245     for( fsarg = ( void * ) header + sizeof( *header );
246         fsarg->type != FSE_ARG_DONE;
247         fsarg = ( void * ) fsarg + ( sizeof( uint16_t ) * 2 + fsarg->length ) ) {
248
249         // handle header & non-file arguments
250         switch( fsarg->type ) {
251             case FSE_ARG_INT32:
252                 ev->eid = ( int64_t ) fsarg->v.int32;
253                 break;
254             case FSE_ARG_INT64:
255                 ev->eid = fsarg->v.int64;
256                 break;
257             case FSE_ARG_STRING:
258                 // move up to the next file arg
259                 if( file == -1 && ( ev->event & SFE_ARG ) ) {
260                     file = 0;
261                     arg = ( struct argsfe * ) ( output + sizeof( struct sfe ) );
262                 }
263                 else if( file == 0 && ( ev->event & SFE_ARGS ) == SFE_ARGS ) {
264                     file = 1;

```

```

265         arg = ( struct argsfe * )
266             ( output + sizeof( struct sfe ) +
267               sizeof( struct argsfe ) + arg->length );
268     }
269     else {
270         file = -2;
271         arg = NULL;
272         break;
273     }
274
275     // copy string, pad to 8-byte boundary
276     arg->length = stpcpy( arg->path, fsarg->v.str ) - arg->path;
277     arg->length = ( ( arg->length >> 3 ) + 1 ) << 3;
278     break;
279 default:
280     // if we don't have a file, ignore specific file args
281     if( file < 0 ) {
282         continue;
283     }
284 }
285
286 // handle specific file args
287 switch( fsarg->type ) {
288 case FSE_ARG_DEV:
289     arg->dev = fsarg->v.dev;
290     break;
291 case FSE_ARG_INO:
292     arg->ino = fsarg->v.inode;
293     break;
294 case FSE_ARG_MODE:
295     arg->mode = fsarg->v.mode;
296     break;
297 case FSE_ARG_UID:
298     arg->uid = fsarg->v.uid;
299     if( ev->uid == -1 ) {
300         ev->uid = arg->uid;
301     }
302     break;
303 case FSE_ARG_GID:
304     arg->gid = fsarg->v.gid;
305     if( ev->gid == -1 ) {
306         ev->gid = arg->gid;
307     }
308     break;
309 }
310 }
311
312 // fix uid/gid if we couldn't use procinfo,
313 // or didn't use one from an argument (rare)
314 if( ev->uid == -1 ) {
315     ev->uid = 0;
316 }
317 if( ev->gid == -1 ) {
318     ev->gid = 0;
319 }
320
321 // return total consumption
322 return ( void * ) fsarg - ( void * ) header + sizeof( uint16_t );
323
324 }
325
326

```

```

327 /**
328  * Translates a given fsevent type to the corresponding sfe type.
329  */
330 uint32_t sfe_type( int32_t fsevent ) {
331     switch( fsevent ) {
332     case FSE_CREATE_FILE:         return SFE_CREATE;
333     case FSE_DELETE:             return SFE_DELETE;
334     case FSE_STAT_CHANGED:       return SFE_STAT;
335     case FSE_RENAME:             return SFE_RENAME;
336     case FSE_CONTENT_MODIFIED:   return SFE_MODIFIED;
337     case FSE_EXCHANGE:          return SFE_EXCHANGE;
338     case FSE_FINDER_INFO_CHANGED: return SFE_FINDER;
339     case FSE_CREATE_DIR:        return SFE_MKDIR;
340     case FSE_CHOWN:             return SFE_CHOWN;
341     }
342     return -1;
343 }
344
345 /**
346  * Register volume information for the volume mounted at the given path.
347  */
348 void register_volume_path( const char *path ) {
349
350     int fd = -1;
351     struct stat info;
352     CFUOIDRef uuid;
353     CFUOIDBytes bytes;
354     uint8_t output[SFE_MAXLEN];
355     struct eyefs request;
356     struct sfe *outputp = ( struct sfe * ) output;
357     struct argsfe *arg = ( struct argsfe * ) ( output + sizeof( struct sfe ) );
358
359     // try to open the device as a fd (lock it from being unmounted)
360     if( -1 == ( fd = open( path, O_RDONLY ) ) ) {
361         fprintf( stderr, "eyemgr: couldn't open device at path '%s'\n", path );
362         return;
363     }
364
365     // try to grab the device node
366     if( fstat( fd, &info ) ) {
367         fprintf( stderr, "eyemgr: couldn't stat %s\n", path );
368         goto fallout;
369     }
370     fprintf( stderr, "eyemgr: register_volume_path got '%s' (dev %d)\n", path, info.st_dev );
371
372     // set up request
373     request.count = 1;
374     request.version = EYEFS_VERSION;
375     request.store.ptr = ( void * ) output;
376
377     // clear output struct (probably not required)
378     memset( ( void * ) output, 0, SFE_MAXLEN );
379     outputp->event = SFE_MOUNT;
380     arg->dev = info.st_dev;
381     request.size = SFE_SIZE( outputp );
382
383     // try to grab the attr
384     if( -1 == fgetxattr( fd, EYEFS_TLD ".", EYEFS_UUID, ( void * ) &arg->uuid, 16, 0, 0 ) ) {
385         // assume ENOATTR
386
387         if( errno != ENOATTR ) {
388             fprintf( stderr, "eyemgr: non-ENOATTR errno (%d) on getting uuid, path = %s\n",

```

```

389         errno, EYEFS_TLD "." EYEFS_UUID );
390     }
391
392     // copy a ref to the device's FSEvents uuid
393     if( NULL == ( uuid = FSEventsCopyUUIDForDevice( info.st_dev ) ) ) {
394         uuid = CFUUIDCreate( NULL );
395     }
396
397     bytes = CFUUIDGetUUIDBytes( uuid );
398     memcpy( arg->uuid, ( char * ) &bytes, 16 );
399     printf( "eyemgr: got id, %llu.%llu\n",
400            *( ( ( uint64_t * ) &arg->uuid ) + 1 ), *( ( uint64_t * ) &arg->uuid ) );
401     CFRelease( uuid );
402
403     // alternatively, generate a random uuid
404     /*
405     else {
406         int i;
407         for( i = 0; i < sizeof( arg->uuid ); i += sizeof( int ) ) {
408             *( ( int * ) ( ( ( char * ) arg->uuid ) + i ) ) = rand();
409         }
410         printf( "eyemgr: generated random id, %llu.%llu\n",
411                *( ( ( uint64_t * ) &arg->uuid ) + 1 ), *( ( uint64_t * ) &arg->uuid ) );
412     }
413     */
414
415     // set xattr
416     if( -1 == fsetxattr( fd,
417                        EYEFS_TLD "." EYEFS_UUID,
418                        ( void * ) &arg->uuid, 16, 0, XATTR_CREATE ) ) {
419         fprintf( stderr, "eyemgr: couldn't set uuid on new volume,"
420                " errno=%d\n", errno );
421         goto fallout;
422     }
423
424     }
425     else {
426         printf( "eyemgr: found id on disk, %llu.%llu\n",
427                *( ( ( uint64_t * ) &arg->uuid ) + 1 ), *( ( uint64_t * ) &arg->uuid ) );
428     }
429
430     // try to push mount event
431     if( ioctl( eyefs, EYEFS_PUSH, &request ) ) {
432         fprintf( stderr, "eyemgr: register_volume_path could not ioctl() with mount info, eyefs=%d, errno=%d\n", eyefs, errno );
433     }
434
435     fallout:
436     if( fd > -1 ) {
437         close( fd );
438     }
439 }
440
441 /**
442  * Handle FSEvents callbacks, and filter out just the mounts and unmounts.
443  */
444 void fsevents_callback(
445     ConstFSEventStreamRef streamRef,
446     void *clientCallBackInfo,
447     size_t numEvents,
448     void *eventPaths,
449     const FSEventStreamEventFlags eventFlags[],

```

```

451     const FSEventStreamEventId eventIds[] ) {
452
453     int i;
454     char **paths = eventPaths;
455
456     for( i = 0; i < numEvents; ++i ) {
457
458         // handle a mount event
459         if( eventFlags[i] & kFSEventStreamEventFlagMount ) {
460             register_volume_path( paths[i] );
461         }
462
463         // we don't do anything with this at the moment/
464         // - can we, since we can't get the device node now?
465         else if( eventFlags[i] & kFSEventStreamEventFlagUnmount ) {
466             printf( "UNMOUNT: %s\n", paths[i] );
467         }
468     }
469 }
470
471 /**
472  * Handles all mounts/unmounts.
473  */
474 void *mount_mgr( __unused void *x ) {
475
476     CFStringRef root = CFSTR( "/" );
477     CFArrayRef pathsToWatch = CFArrayCreate( NULL, ( const void ** ) &root, 1, NULL );
478     FSEventStreamRef stream;
479     CFAbsoluteTime latency = 0.2;
480     Boolean ret;
481     int nDevices, index;
482     struct statfs *fsList;
483
484     // set up FSEvents stream
485     stream = FSEventStreamCreate( NULL,
486                                  &fsevents_callback, NULL, pathsToWatch,
487                                  kFSEventStreamEventIdSinceNow, latency,
488                                  kFSEventStreamCreateFlagNoDefer );
489     FSEventStreamScheduleWithRunLoop(
490         stream, CFRunLoopGetCurrent(), kCFRunLoopDefaultMode );
491     ret = FSEventStreamStart( stream );
492     if( !ret ) {
493         fprintf( stderr, "eyefs=%d, errno=%d\n", eyefs, errno );
494     }
495     // handle file systems that already existed
496     if( ( nDevices = getfsstat( NULL, 0, MNT_NOWAIT ) ) < 0 ) {
497         fprintf( stderr, "fatal; getfsstat failed\n" );
498         exit( 1 );
499     }
500
501     // allocate space for the fslist
502     if( NULL == ( fsList = ( struct statfs * )
503                malloc( nDevices * sizeof( struct statfs ) ) ) ) {
504         perror( "malloc" );
505         exit( 2 );
506     }
507
508     // obtain information about existing file systems
509     if( ( nDevices = getfsstat( fsList, nDevices * sizeof( struct statfs ), MNT_NOWAIT ) ) < 0 ) {
510         fprintf( stderr, "fatal; getfsstat failed with full request\n" );
511         exit( 3 );
512     }

```

```

513     for( index = 0; index < nDevices; ++index ) {
514         register_volume_path( fsList[index].f_mntonname );
515     }
516     free( fsList );
517
518     // run loop
519     printf( "eyemgr: mount_mgr setup ok (%d)\n", ret );
520     CFRRunLoopRun();
521
522     // shutdown!
523     CFRRelease( pathsToWatch );
524     CFRRelease( stream );
525
526     return 0;
527
528 }

```

fastlock.c

```

1
2 #include "fastlock.h"
3
4 /**
5  * Obtain a shared lock on the integer at the given location. Returns the
6  * positive lock count, or -1 if locked exclusively.
7  */
8 int FastSIntLock( fastlock_t *v ) {
9     int32_t value;
10    while( ( value = *v ) >= 0 ) {
11        if( OSCompareAndSwap( value, value + 1, ( UInt32 * ) v ) ) {
12            return value + 1;
13        }
14    }
15    return -1;
16 }
17
18 /**
19  * Releases a shared lock on the integer at the given location. Returns the
20  * positive (or zero) lock count, or -1 if locked exclusively (or unlocked).
21  */
22 int FastSIntUnlock( fastlock_t *v ) {
23     int32_t value;
24     while( ( value = *v ) > 0 ) {
25         if( OSCompareAndSwap( value, value - 1, ( UInt32 * ) v ) ) {
26             return value - 1;
27         }
28     }
29     return -1;
30 }
31
32 /**
33  * Obtains an exclusive lock with the given integer location. Returns true if
34  * successful, false otherwise.
35  */
36 Boolean FastSIntLockEx( fastlock_t *v ) {
37     return OSCompareAndSwap( 0, -1, ( UInt32 * ) v );
38 }

```

```

39
40 /**
41  * Releases an exclusive lock with the given integer location. returns true if
42  * successful, false otherwise.
43  */
44 Boolean FastSIntUnlockEx( fastlock_t *v ) {
45     return OSCompareAndSwap( -1, 0, ( UInt32 * ) v );
46 }

```

fastlock.h

```

1
2 #include "fastlock.h"
3
4 /**
5  * Obtain a shared lock on the integer at the given location. Returns the
6  * positive lock count, or -1 if locked exclusively.
7  */
8 int FastSIntLock( fastlock_t *v ) {
9     int32_t value;
10    while( ( value = *v ) >= 0 ) {
11        if( OSCompareAndSwap( value, value + 1, ( UInt32 * ) v ) ) {
12            return value + 1;
13        }
14    }
15    return -1;
16 }
17
18 /**
19  * Releases a shared lock on the integer at the given location. Returns the
20  * positive (or zero) lock count, or -1 if locked exclusively (or unlocked).
21  */
22 int FastSIntUnlock( fastlock_t *v ) {
23     int32_t value;
24     while( ( value = *v ) > 0 ) {
25         if( OSCompareAndSwap( value, value - 1, ( UInt32 * ) v ) ) {
26             return value - 1;
27         }
28     }
29     return -1;
30 }
31
32 /**
33  * Obtains an exclusive lock with the given integer location. Returns true if
34  * successful, false otherwise.
35  */
36 Boolean FastSIntLockEx( fastlock_t *v ) {
37     return OSCompareAndSwap( 0, -1, ( UInt32 * ) v );
38 }
39
40 /**
41  * Releases an exclusive lock with the given integer location. returns true if
42  * successful, false otherwise.
43  */
44 Boolean FastSIntUnlockEx( fastlock_t *v ) {
45     return OSCompareAndSwap( -1, 0, ( UInt32 * ) v );
46 }

```


Appendix B

Lighthouse

The code below is the complete implementation of Lighthouse in Python.

core.py

This describes the core of Lighthouse, and mirrors the defined algorithm most explicitly. It does not provide interactive functionality, however.

```
1 #!/usr/bin/env python
2
3 from __future__ import with_statement
4 import eyedropper, schema, re, os, time, uuid, stat, threading, sys, signal
5 import logging
6 from util import *
7
8 logging.basicConfig()
9
10 class Lighthouse( object ):
11     """ Lighthouse watches your files. To a big extent, it is the middle-man
12     between direct disk access and behaviour (lighthouse/inodes), and the
13     node-based database structure which conceptualises nodes and their
14     relationships. """
15
16     def __init__( self ):
17         """ Configures a new Lighthouse object. """
18         self.devices = DeviceSet()
19
20     def watch( self ):
21         """ A helper method to watch for events until a KeyboardInterrupt.
22         This method doesn't maintain a state over the last seen event, so
23         calling it multiple times will mean events are dropped. """
24         try:
25             for event in eyedropper.listen():
26                 self.act( event )
27         except KeyboardInterrupt:
28             pass
29
30     def threaded( self ):
31         """ Waits for events to be pushed into our queue, then acts on them. """
32
33         # set up shared objects
34         condition = threading.Condition()
35
36         # threaded function
37         def func():
38             while func.active:
39                 with condition:
40                     while func.active and not len( func.queue ):
41                         condition.wait()
42                     mine = func.queue
43                     func.queue = []
44                     for i in mine:
45                         self.act( i )
46         func.queue = []
47         func.active = True
48
49         # run function, pass events through
50         th = threading.Thread( target = func )
51         th.start()
52         try:
53             for event in eyedropper.listen():
54                 assert event is not None
55                 with condition:
56                     func.queue.append( event )
57                     condition.notify()
```

```
58
59     # fallout conditions
60     except KeyboardInterrupt:
61         print "## Ctrl-C."
62     except AssertionError:
63         pass
64
65     # shutdown
66     func.active = False
67     with condition:
68         condition.notify()
69     th.join()
70     print "## Shutdown OK."
71
72 def act( self, event ):
73     """ Act on a single event. """
74
75     # ignore our own pid.
76     if event.pid == os.getpid():
77         pass
78
79     # standard unix rename - if target has an inode, it 'exists'
80     elif event.type == eyedropper.RENAME:
81
82         # this file exists, and will be baledeted!
83         if event.target.inode:
84             self.release( event.target )
85
86         # seen the original, but with a new name
87         node = self.seen( event.focus, path = event.target.path )
88
89     # content flip
90     elif event.type == eyedropper.EXCHANGE:
91         # TODO: remember how this works
92         pass
93
94     # otherwise, just observe (or release) the file
95     else:
96         if event.type == eyedropper.DELETE:
97             self.release( event.focus )
98         else:
99             node = self.seen( event.focus )
100
101 def release( self, file ):
102     """ Handles notifications that files have been released. Does not
103     attempt to interact with the file on-disk, as it has been released,
104     and we might interact with a new file in the same place. """
105
106     if file.uid != os.getuid():
107         return
108
109     # grab device structure
110     device = self.devices.get( file.uuid, file.device, file.path )
111     if not device:
112         return None
113     assert device.uuid is not None
114
115     # try to find a matching node
116     nodes = schema.Node.selectBy( inode = file.inode, device = device.uuid )
117     try:
118         node = max( nodes, key = lambda f: f.seen )
119
```

```

120     except ValueError:
121         return None
122
123     node.apply( schema.Destroyed() )
124     print "-- Releasing node (last seen at) '%s' (%s)"
125           % ( os.path.join( node.dirname, node.basename ), node.uuid )
126
127     # return node itself
128     return node
129
130 def props( self, path ):
131     """ Returns a small tuple of properties of the file at the given
132     path. Format is ( created, modified, size ). """
133
134     # standard args
135     result = os.stat( path )
136     modified = result.st_mtime
137     size = result.st_size
138     created = Created( path )
139
140     # modify directory results /slightly/ so we don't have so many changes
141     if stat_S_ISDIR( result.st_mode ):
142         modified = created
143         size = -1
144
145     # return tuple
146     return( created, modified, size )
147
148 def brand( self, node, force = False ):
149     """ Brand a given LocalDevice node with an ID, or grab the current
150     branding. """
151
152     try:
153         if force:
154             raise AttributeError
155
156         id = node.id
157
158         # prevent sqlalchemy from breaking
159         if len( filter( lambda f: ord( f ) > 127, id ) ):
160             id = None
161             raise AttributeError
162
163         # other weird stuff
164         if id is None:
165             raise AttributeError
166
167     except AttributeError:
168         try:
169             #
170             id = uuid.uuid1().bytes
171             id = uuid.uuid1().hex
172             node.id = id
173         except AttributeError:
174             raise OSError
175
176     return id
177
178 def seen( self, file, path = None ):
179     """ Handles notifications that files have been observed. From a
180     high-level point of view, this method doesn't care about how this
181     has occurred. Returns the found (current) node. """

```

```

182     stime = time.time()
183
184     if file.uuid != os.getuid():
185         return
186
187     # fix presented path, grab relevant device struct
188     if path is None:
189         path = file.path
190     device = self.devices.get( file.uuid, file.device, path )
191     if not device:
192         return None
193     assert device.uuid is not None
194
195     # how can we access this file?
196     if device.volfs:
197         apath = "/.vol/%d/%d" % ( file.device, file.inode )
198     else:
199         apath = path
200
201     # grab parent data
202     try:
203         parent = os.stat( os.path.dirname( path ) )
204         if parent.st_dev == file.device and parent.st_ino != file.inode:
205             parent = parent.st_ino
206         else:
207             parent = 0 # we are the top level node
208     except ( OSError, IOError ):
209         parent = None # parent is undefined
210
211     # get information about the file in question (if we OSError, the file has gone)
212     try:
213         created, modified, size = self.props( apath )
214         dirname = unicode( os.path.dirname( path ), 'utf-8' )
215         basename = unicode( os.path.basename( path ), 'utf-8' )
216         id = self.brand( device[file.inode] )
217     except ( OSError, IOError ):
218         print "-- Node disappeared '%s' (?)" % ( path, )
219         return
220
221     # create and compare all immutable evidence on this node
222     try:
223         node = schema.Node.find( id )
224         assert node.inode is not None
225         assert node.inode == file.inode
226         assert node.pinode == parent
227         assert node.device == device.uuid
228         assert node.size == size
229         assert node.modified == modified
230         assert node.created == created
231         #
232         assert node.dirname == dirname
233         assert node.basename == basename
234         print "-- Touched node '%s' (%s)" % ( path, id )
235
236     # update our node structure
237     except AssertionError:
238
239         logging.exception( "Node Changed (%r, %r)", node.basename, basename )
240
241     # something changed - create a new derived node
242     dbtime = time.time()
243     if node.inode:
244         id = self.brand( device[file.inode], force = True )

```

```

244         replacement = schema.Node( uuid = id )
245
246     # apply relevant evidence
247     if node.inode != file.inode or node.device != device.uuid:
248         replacement.apply(
249             schema.CopyCreation( parent = node )
250         )
251         print "-- Creating copied node '%s' (%s)"
252             % ( path, id )
253     else:
254         replacement.apply(
255             schema.DerivedFrom( parent = node )
256         )
257         node.apply( schema.Destroyed() )
258         print "-- Creating derived/updated node '%s' (%s)"
259             % ( path, id )
260
261     # focus on our new node
262     node = replacement
263
264     else:
265         print "-- Creating completely new node '%s' (%s)"
266             % ( path, id )
267
268     # add same location hints
269     nodes = schema.Node.selectBy(
270         pinode = parent, device = device.uuid, basename = basename
271     )
272     try:
273         target = max( nodes, key = lambda f: f.seen )
274         node.apply( schema.SameLocation( previous = target ) )
275         print "-- Added SameLocation evidence directed at previous "
276             + "node (%s)" % ( target.uuid, )
277     except ValueError:
278         pass
279
280     # configure the node
281     node.inode = file.inode
282     node.pinode = parent
283     node.device = device.uuid
284     node.size = size
285     node.modified = int( modified )
286     node.created = int( created )
287     node.dirname = dirname
288     node.basename = basename
289     print "dbtime took %f" % ( time.time() - dbtime )
290
291     print "took %f" % ( time.time() - stime )
292
293     # return the node
294     return node

```

schema.py

The SQLAlchemy ?? database schema, providing *nodes* and *evidence*.

```

1 #!/usr/bin/env python
2
3 """ This file defines the basic set of evidence available to nodes within
4 the Lighthouse/Family system (aka Sam's thesis). Evidence may be attached

```

```

5 to any particular Node instance within the system, and is based around the
6 core idea that it defines why this node exists. """
7
8 from sqlalchemy import *
9 from sqlalchemy.inheritance import *
10 import uuid, os
11
12 class Node( SQLAlchemy ):
13     """ Each node represents an immutable file within our graph. A single
14     node may be dead or alive; its reasoning for being created or destroyed
15     is however attached to its micro-reasons. Since they are immutable, a
16     node may exist for current and prior device/inode combinations. """
17
18     # hex uuid, timestamp, evidence link
19     uuid = StringCol( length = 32, alternateID = True )
20     seen = DateTimeCol( default = DateTimeCol.now, notNone = True )
21     evidence = MultipleJoin( 'BaseEvidence' )
22
23     # device/fs attributes and timestamps
24     inode = IntCol( default = None )
25     device = StringCol( length = 32, default = None )
26     pinode = IntCol( default = None )
27     size = IntCol( default = None )
28     modified = IntCol( default = None )
29     created = IntCol( default = None )
30
31     # path attributes
32     dirname = UnicodeCol( default = None )
33     basename = UnicodeCol( default = None )
34
35     # indexes
36     inodeIndex = DatabaseIndex( 'inode' )
37     inodeDeviceIndex = DatabaseIndex( 'inode', 'device' )
38     parentDeviceName = DatabaseIndex( 'pinode', 'device', 'basename' )
39
40     # get a node for this uuid
41     @staticmethod
42     def find( id ):
43         """ Get a node for the passed UUID, no matter what the cost. """
44         try:
45             node = Node.byUuid( id )
46         except SQLAlchemyObjectNotFound:
47             node = Node( uuid = id )
48         return node
49
50     # apply given evidence to this object
51     def apply( self, evidence, *args ):
52         for v in [ evidence ] + list( args ):
53             assert isinstance( v, BaseEvidence )
54             v.node = self
55
56     def __str__( self ):
57         return self.uuid
58
59     def __repr__( self ):
60         return( "<Node uuid=\"%s\" seen=%s destroyed=%s device=\"%s\" inode=%s pinode=%s "
61             + "size=%s modified=%s created=%s path=\"%s\">" ) % (
62             self, self.seen, self.destroyed, self.device, self.inode, self.pinode,
63             self.size, self.modified, self.created, os.path.join( self.dirname, self.basename )
64         )
65
66     @property

```

```

67 def ancestors( self ):
68     """ Return all direct ancestors of this node. """
69
70     output = set()
71
72     temp = list( DerivedFrom.select( BaseEvidence.q.node == self ) )
73     temp = map( lambda f: f.parent, temp )
74     output.update( temp )
75
76     temp = list( SameLocation.select( BaseEvidence.q.node == self ) )
77     temp = map( lambda f: f.previous, temp )
78     output.update( temp )
79
80     # success!
81     return output
82
83 @property
84 def descendants( self ):
85     """ Return all direct descendants of this node. """
86     # generate list of direct children (evidence)
87     output = set()
88     output.update(
89         list( DerivedFrom.select( DerivedFrom.q.parent == self ) )
90     )
91     output.update(
92         list( SameLocation.select( SameLocation.q.previous == self ) )
93     )
94
95     # map this to actual nodes
96     output = set( map( lambda f: f.node, output ) )
97
98     # success!
99     return output
100
101 @property
102 def destroyed( self ):
103     """ Return the time at which this node was destroyed (or False). """
104     evidence = list( Destroyed.select( BaseEvidence.q.node == self ) )
105     if len( evidence ):
106         return evidence[0].timestamp
107     return False
108
109 class BaseEvidence( InheritableSQLObject ):
110     """ The root evidence object. Can be instantiated on its own, but would
111     mean extremely little. """
112
113     # node target, timestamp
114     node = ForeignKey( 'Node', default = None )
115     timestamp = DateTimeCol( default = DateTimeCol.now, notNone = True )
116
117     # indexes
118     nodeIndex = DatabaseIndex( 'node' )
119
120 class SameLocation( BaseEvidence ):
121     """ This file is at the same location as a previous file. """
122     previous = ForeignKey( 'Node', notNone = True )
123
124 class Replaced( SameLocation ):
125     """ This file replaced a previous node, causing it to be destroyed. """
126     pass
127
128 class Destroyed( BaseEvidence ):

```

```

129     """ This node has been observed as destroyed. """
130     pass
131
132 class DerivedFrom( BaseEvidence ):
133     """ This node was directly derived (modified) from a previous node. """
134     parent = ForeignKey( 'Node', notNone = True )
135
136 class CopyCreation( DerivedFrom ):
137     """ This node is a direct copy of a previous node. """
138     pass

```

interactive.py

Sets up the local Lighthouse database using SQLAlchemy. Included by other, interactive tools.

```

1 #!/usr/bin/env python
2
3 import sqlalchemy, core, schema, os
4
5 # define connection
6 path = os.path.expanduser( "~/.lighthouse.db" )
7 sqlalchemy.sqlhub.processConnection = sqlalchemy.connectionForURI( "sqlite://" + path )
8 #sqlalchemy.sqlhub.processConnection.debug = True
9
10 # try to create tables
11 tables = [ 'Node', 'BaseEvidence', 'SameLocation', 'Replaced', 'Destroyed', 'DerivedFrom', 'CopyCreation' ]
12 for i in tables:
13     try:
14         schema.__dict__[i].createTable()
15         print "Created table '%s': ok" % i
16     except:
17         if __name__ == "__main__":
18             print "Created table '%s': already exists" % i
19

```

test.py

This sets up and then runs the core of Lighthouse; invoke this interactive tool to begin monitoring events.

```

1 #!/usr/bin/env python
2
3 import sqlalchemy, core, schema, os
4 import interactive
5
6 if __name__ == "__main__":
7
8     # watch for events
9     house = core.Lighthouse()
10     try:
11         house.threaded()
12     except KeyboardInterrupt:
13         pass
14

```

related.py

Provides the main testing framework as used within our evaluation, and described in our implementation. An interactive tool.

```
1 #!/usr/bin/env python
2
3 import sys, os, interactive, schema, util
4
5 # setups
6 assert __name__ == "__main__"
7 assert len( sys.argv ) == 2
8 target = sys.argv[1]
9
10 # grab inode/device
11 s = os.stat( target )
12 inode, device = s.st_ino, s.st_dev
13
14 # grab database node
15 node = schema.Node.byUuid( util.Node( device, inode ).id )
16
17 # grab all 'family' nodes
18 nodes = [ node ]
19 i = 0
20 while i < len( nodes ):
21     node = nodes[i]
22
23     for e in node.ancestors:
24         if e not in nodes:
25             nodes.append( e )
26
27     for e in node.descendants:
28         if e not in nodes:
29             nodes.append( e )
30
31     i += 1
32
33 # display all
34 print "got %d results" % len( nodes )
35 print
36 for i in nodes:
37     print repr( i )
38     print "ancestors=%d descendants=%d" % ( len( i.ancestors ), len( i.descendants ) )
39     print
```

util.py

Mostly provides code used to help identify files over time, as discussed in our implementation. Used internally by other parts of the code.

```
1 #!/usr/bin/env python
2
3 import warnings
4 warnings.simplefilter( "ignore", DeprecationWarning )
5
6 import os, pickle, re
7 import uuid
8 from uuid import UUID
9 from xattr import xattr
```

```
10 import MacOS, macfs
11
12 class DeviceSet( object ):
13     """ A wrapper object to help present all devices on a system. """
14     # FIXME: obviously this doesn't handle inserts/removal of devices at the moment.
15
16     def __init__( self ):
17         self.devices = {}
18
19     def get( self, uuid, device, path ):
20         """ Get a device, directly based on uuid, device no, and reference path. """
21
22         # hackery laem
23         real = long()
24         for i in uuid:
25             real <<= 8
26             real += ord( i )
27         uuid = hex( long( real ) )[2:-1]
28         del real
29
30         if not uuid:
31             return None
32
33         if device not in self.devices or self.devices[device].uuid != uuid:
34             try:
35                 self.devices[device] = LocalDevice( uuid, device, path )
36                 print "new device created: dev %d, uuid %s, ref path = %s" % ( device, uuid, path )
37             except OSError:
38                 return None
39
40         return self.devices[device]
41
42     def __contains__( self, device ):
43         """ Does this set contain the given device node? """
44         return device in self.devices
45
46 class LocalDevice( object ):
47     """ A device object which represents a local, accessible device. Allows
48     users to interact with paths or inodes. """
49
50     def __init__( self, uuid, device, path ):
51         """ Create a new LocalDevice with reference arguments """
52
53         if not uuid:
54             raise OSError( "Blank UUID passed." )
55
56         self.uuid = uuid
57         self.device = os.stat( path ).st_dev
58
59         if self.device != device:
60             raise OSError( "Path/device mismatch." )
61
62         try:
63             self.prop = self.props( "/.vol/%d/%" % ( self.device, ) )
64             self.volfs = True
65         except OSError:
66             self.volfs = False
67         if not path:
68             raise Exception( "This device does not exist or support volfs,"
69                             + " so a path must be specified." )
70         self.prop = self.props( path )
71         self.mount = self.prop['Mount Point']
```

```

72
73
74 def __contains__( self, inode ):
75     """ Does this device currently contain the passed inode? """
76
77     # throw out if volfs is not supported
78     if not self.volfs:
79         raise Exception( "volfs is not supported on this device!" )
80
81     # throw out on badly typed submissions
82     if not isinstance( inode, ( int, long ) ) or inode < 0:
83         raise TypeError( "A LocalDevice only supports positive "
84             + "integer inodes." )
85
86     # fix inode request
87     if not inode:
88         inode = "0"
89
90     # return the status of this request
91     return os.path.exists( "/.vol/%d/%s" % ( self.device, inode ) )
92
93 def __getitem__( self, inode ):
94     """ Returns an inode structure from the LocalDevice. """
95
96     # inode does not exist
97     if inode not in self:
98         return None
99
100     # fix inode request
101     if not inode:
102         inode = "0"
103
104     # return a relevant node structure
105     return Node( self.device, inode )
106
107 @staticmethod
108 def props( path ):
109     """ Returns properties of the device mounted at the given path. """
110     output = {}
111     source = os.popen( "diskutil info %s" % path ).read().split( "\n" )
112     for i in source:
113         try:
114             k, v = i.split( ":", 1 )
115             output[k.strip()] = v.strip()
116         except:
117             pass
118     if "Usage" in output:
119         raise OSError( "No such file or directory: '%s'" % ( path, ) )
120     return output
121
122 class Node( object ):
123     """ A simple xattr wrapper that should only be used by LocalDevice. """
124     _FORMAT = "com.moofco.lighthouse.%s"
125
126     _xattr = None
127
128 def __init__( self, device, inode ):
129     self._xattr = xattr( "/.vol/%d/%s" % ( device, inode ) )
130
131 def __getattr__( self, k ):
132     if k.startswith( "_" ):
133         return super( Node, self ).__getattr__( k )
134     path = Node._FORMAT % k
135     if path not in self._xattr:
136         raise AttributeError, k
137     try:
138         return pickle.loads( self._xattr[path] )
139     except:
140         raise AttributeError, k
141
142 def __setattr__( self, k, v ):
143     if k.startswith( "_" ):
144         super( Node, self ).__setattr__( k, v )
145     else:
146         path = Node._FORMAT % k
147         if v is None:
148             del self._xattr[path]
149         else:
150             self._xattr[path] = pickle.dumps( v )
151
152 def __delattr__( self, k ):
153     if k.startswith( "_" ):
154         super( Node, self ).__delattr__( k, v )
155     else:
156         self.__setattr__( k, None )
157
158 def Created( path ):
159     """ Return the creation time of a specified path. Mac specific. """
160     # weird base case that tends to fail out
161     if re.search( "~\\.vol/(\\d+)/0$", path ):
162         cmd = "diskutil info %s | grep '^ Mount Point'" % path
163         output = os.popen( cmd ).read()
164         try:
165             path = re.search( "^ Mount Point:\\s+(\\S.*)$", output ).group( 1 )
166         except:
167             raise Exception( "Can't find mount point of root inode: '%s'" % path )
168
169     # normal case (easy)
170     try:
171         return float( macfs.FSSpec( path ).GetDates()[0] )
172     except MacOS.Error:
173         raise OSError( "No such file or directory: '%s'" % path )
174
175 # export list
176 __all__ = [ "DeviceSet", "LocalDevice", "Created" ]

```