

# Formally Specifying POSIX File Systems

**Gian Ntzik**, Pedro da Rocha Pinto and Philippa Gardner

Imperial College London

July 16, 2015

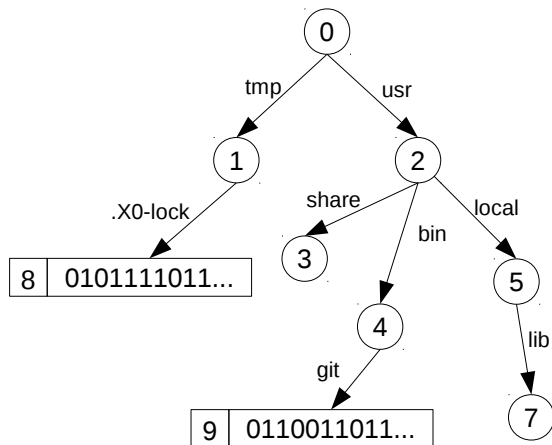
# POSIX File Systems

- ▶ POSIX: Portable Operating System Interface
- ▶ Large part devoted to the file system
- ▶ English Specification
  - ▶ Underspecified
  - ▶ Ambiguous
  - ▶ Absence of proper memory model

# Objectives

- ▶ Find formalism suitable for specifying POSIX file system operations
  - ▶ Suitable for clients and implementations
  - ▶ Client reasoning in a program logic
- ▶ Focus a core fragment:
  - ▶ Structural operations: `mkdir`, `rmdir`, `link`, `unlink`, `rename`, ...
  - ▶ IO: `open`, `read`, `write`, `lseek`, `close`, ...

# File System Example

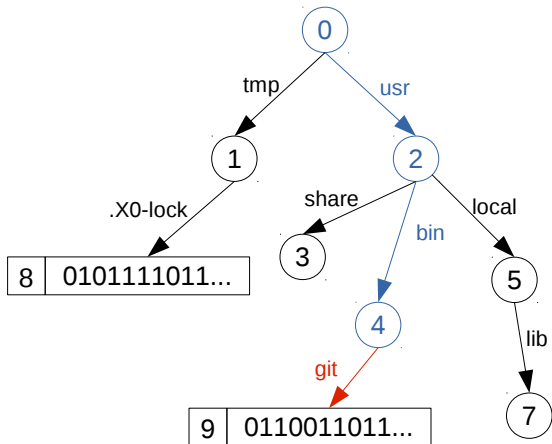


## First Challenge: Atomicity

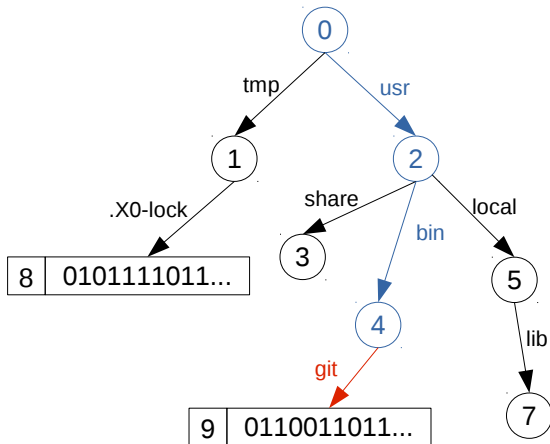
## Example: `unlink`

- ▶ `unlink(p)` : Atomically remove the file identified by path `p`.
- ▶ “Atomic” has an unusual meaning in POSIX.

`unlink(/usr/bin/git)`



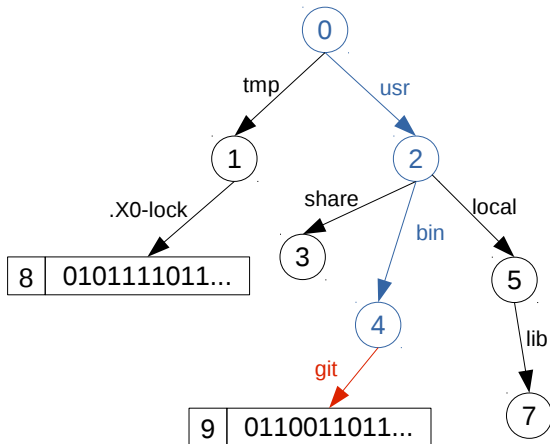
`unlink(/usr/bin/git)`



- ▶ Only removing `git` is required to be atomic.

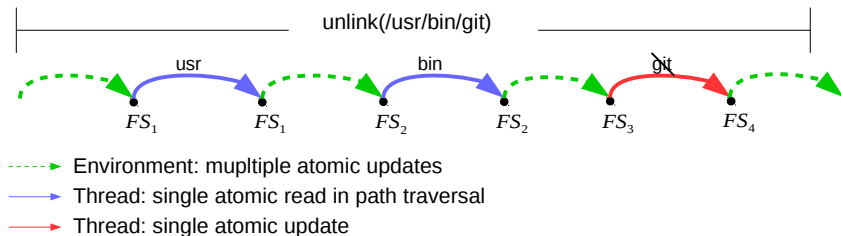


`unlink(/usr/bin/git)`



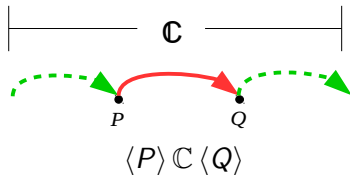
- ▶ Only removing `git` is required to be atomic.
- ▶ Path resolution: a sequence of atomic reads.

# Sequence of atomic actions



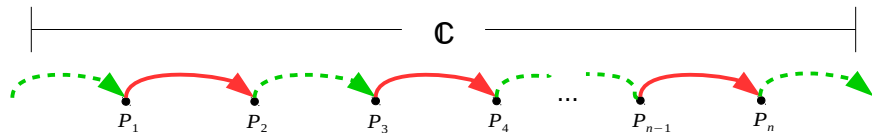
# Single atomic step specification: Atomic Hoare Triple

From TaDA, IRIS, ... we know how to specify a single atomic step.



# Multi-atomic specifications

We extend to multiple atomic steps



$$\mathbb{C} \sqsubseteq \langle (P_1, P_2); (P_3, P_4); \dots; (P_{n-1}, P_n) \rangle$$

# Multi-atomic Program Logic

- ▶ Introduce multi-atomics:

$$\mathbb{C} \sqsubseteq \langle (P_1, P_2); (P_3, P_4); \dots; (P_{n-1}, P_n) \rangle$$

- ▶ Justified by an encoding in IRIS
- ▶ Extend reasoning rules for single atomic steps to multiple steps

# Single Step Equivalence

SINGLEATOMIC:

$$\mathbb{C} \sqsubseteq \langle (P, Q) \rangle \iff \langle P \rangle \mathbb{C} \langle Q \rangle$$

## Sequence Rule

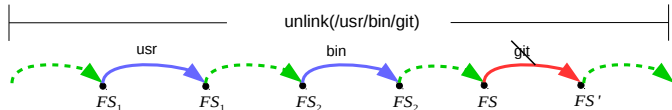
$$\frac{\mathbb{C}_1 \sqsubseteq \langle \dots; (P_1, Q_1) \rangle \quad \mathbb{C}_2 \sqsubseteq \langle (P_2, Q_2); \dots \rangle}{\mathbb{C}_1; \mathbb{C}_2 \sqsubseteq \langle \dots; (P_1, Q_1); (P_2, Q_2); \dots \rangle} \text{MULTI-SEQ}$$

# Stuttering Rule

$$\frac{\mathbb{C} \sqsubseteq \langle \dots; (P, P); (P, Q); \dots \rangle}{\mathbb{C} \sqsubseteq \langle \dots; (P, Q); \dots \rangle} \text{STUTTER}$$

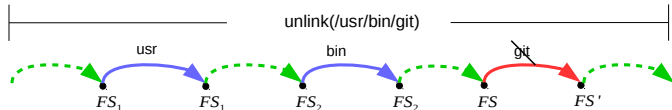


# unlink: Formal Specification



`unlink(/usr/bin/git) ⊆ ⟨resolve(/usr/bin,  $\iota_0, r$ ); ...⟩`

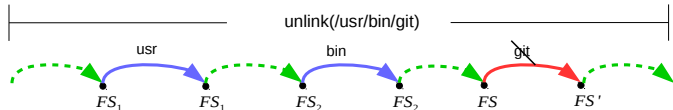
# unlink: Formal Specification



`unlink(/usr/bin/git)  $\sqsubseteq$`

`$\langle$ resolve(/usr/bin,  $\iota_0$ , r); (fs(FS), in(FS, r, git)  $\Rightarrow$  rem(FS, r, git) * ret = 0) $\rangle$`

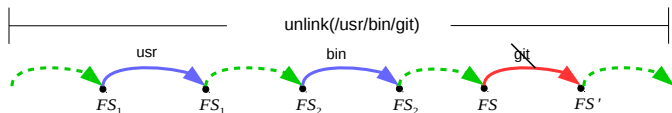
# unlink: Formal Specification



$\text{unlink}(/usr/bin/git) \sqsubseteq$

$\langle \text{resolve}(/usr/bin, \iota_0, r); \left( \text{fs}(FS), \left( \text{in}(FS, r, \text{git}) \Rightarrow \text{rem}(FS, r, \text{git}) * \text{ret} = 0 \right. \right. \right. \\ \left. \left. \left. \wedge \text{out}(FS, r, \text{git}) \Rightarrow \text{fs}(FS) * \text{ret} = -1 * \text{errno} = \text{ENOENT} \right) \right) \rangle$

# unlink: Formal Specification



$\text{unlink}(/usr/bin/git) \sqsubseteq$

$$\left\langle \text{resolve}(/usr/bin, \iota_0, r); \left( \text{fs}(FS), r \in \text{IN} \Rightarrow \left( \text{in}(FS, r, \text{git}) \Rightarrow \text{rem}(FS, r, \text{git}) * \text{ret} = 0 \right. \right. \right. \\ \left. \left. \wedge \text{out}(FS, r, \text{git}) \Rightarrow \text{fs}(FS) * \text{ret} = -1 * \text{errno} = \text{ENOENT} \right) \right. \\ \left. \wedge r \in \text{ERR} \Rightarrow \text{fs}(FS) * \text{ret} = -1 * \text{errno} = r \right. \left. \right) \left. \right\rangle$$

## Second Challenge: Unordered actions

## Specifying unordered actions

- ▶ Example: `rename(p/a, p'/b)`
- ▶ POSIX does not specify in which order `p` and `p'` are resolved

## Specifying unordered actions

- ▶ Example: `rename(p/a, p'/b)`
- ▶ POSIX does not specify in which order `p` and `p'` are resolved
- ▶ We can't do:

$$\text{rename}(p/a, p'/b) \sqsubseteq \langle \text{resolve}(p, \iota_0, r_1); \text{resolve}(p', \iota_0, r_2); \dots \rangle$$

## Specifying unordered actions

- ▶ Example: `rename(p/a, p'/b)`
- ▶ POSIX does not specify in which order `p` and `p'` are resolved
- ▶ Solution:

$$\text{rename}(p/a, p'/b) \sqsubseteq \langle \text{resolve}(p, \iota_0, r_1) \parallel \text{resolve}(p', \iota_0, r_2); \dots \rangle$$



## Parallel Rule

$$\frac{\begin{array}{l} \mathbb{C}_1 \sqsubseteq \langle (P_1, P_2); \dots; (P_{n-1}, P_n) \rangle \\ \mathbb{C}_1 \sqsubseteq \langle (Q_1, Q_2); \dots; (Q_{n-1}, Q_n) \rangle \end{array}}{\mathbb{C}_1 \parallel \mathbb{C}_2 \sqsubseteq \langle (P_1, P_2); \dots; (P_{n-1}, P_n) \parallel (Q_1, Q_2); \dots; (Q_{n-1}, Q_n) \rangle} \text{MULTI-PAR}$$

# Client Applications

- ▶ Lock Files
- ▶ POSIX pipes
- ▶ POSIX Advisory Locks/Record Locking (on-going)
- ▶ Persistent Concurrent Queues (future)

# Lock Files

- ▶ `lock(path)`: atomically create a non-existing lock file at `path`
- ▶ `unlock(path)`: remove the lock file identified by `path`
- ▶ Implemented similarly to spin locks
  - ▶ `open(path, O_CREAT|O_EXCL)` to try to lock
  - ▶ `unlink` to unlock

# Lock File Implementation

```
function lock(path) {  
  do {  
    fd := open(path, O_EXCL|O_CREAT);  
  } while (fd = -1);  
  close(fd);  
}  
  
function unlock(path) {  
  unlink(path);  
}
```

# Heap Based Lock Specification

- ▶ We know how to specify locks on the heap

# Heap Based Lock Specification

- ▶ We know how to specify locks on the heap

$$\{\text{emp}\} \text{makeLock}() \{\text{Lock}(\text{ret}, 0)\}$$
$$\langle \text{Lock}(x, v) \rangle \text{lock}(x) \langle \text{Lock}(x, 1) * v = 0 \rangle$$
$$\langle \text{Lock}(x, 1) \rangle \text{unlock}(x) \langle \text{Lock}(x, 0) \rangle$$

# Heap Based Lock Specification

- ▶ We know how to specify locks on the heap

$$\{\text{emp}\} \text{makeLock}() \{\text{Lock}(\text{ret}, 0)\}$$
$$\langle \text{Lock}(x, v) \rangle \text{lock}(x) \langle \text{Lock}(x, 1) * v = 0 \rangle$$
$$\langle \text{Lock}(x, 1) \rangle \text{unlock}(x) \langle \text{Lock}(x, 0) \rangle$$

- ▶ Ideally, we want the same specification for lock files
- ▶ Using `path` instead of heap address `x`

## Third Challenge: Ownership



## The Heap case

Allocated Heap

makeLock()

Allocated Heap

0

ret

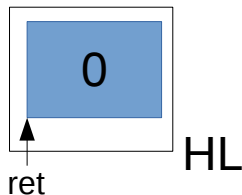
- ▶ The module owns the newly allocated memory

## The Heap case

Allocated Heap

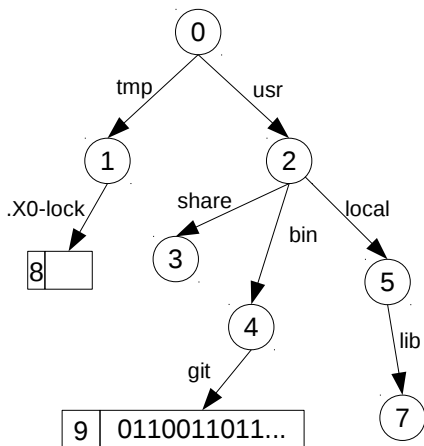
Allocated Heap

makeLock()

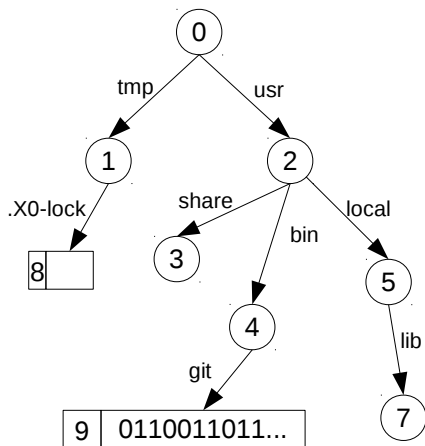


- ▶ The module enforces the sharing protocol

# The File System case

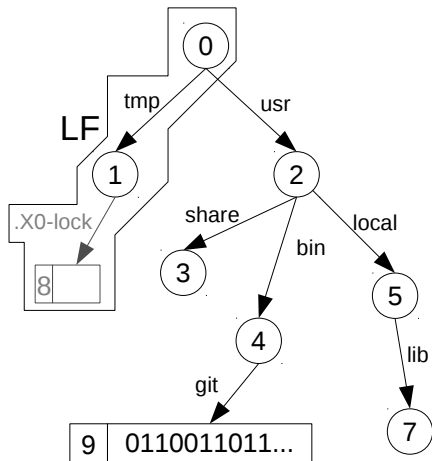


# The File System case

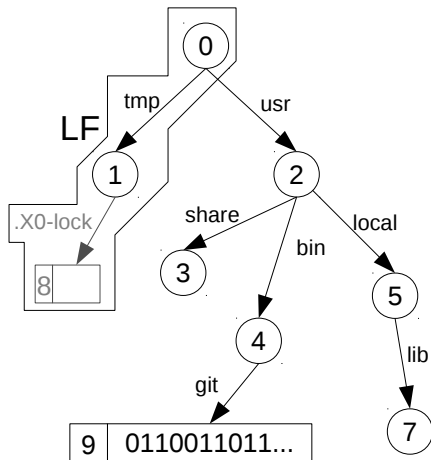


- ▶ No operation to extend with fresh path, unknown to environment
- ▶ Global path address space

# The File System case

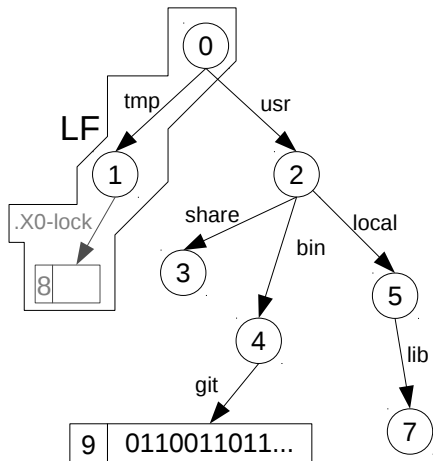


# The File System case



- ▶ Clients must agree on sharing protocol

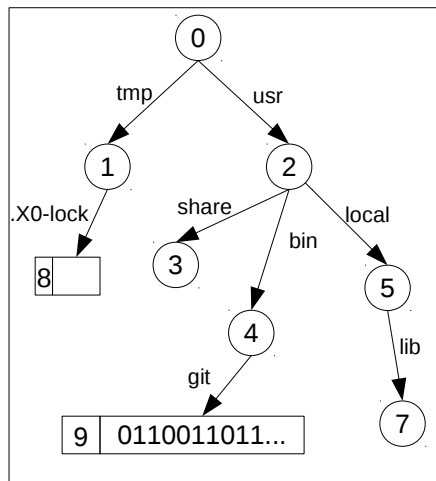
## The File System case



- ▶ Clients must agree on sharing protocol
- ▶ *Cooperative ownership*

# Lock file specification

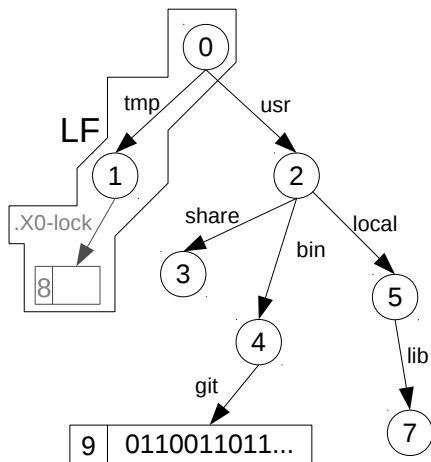
If any client's protocols





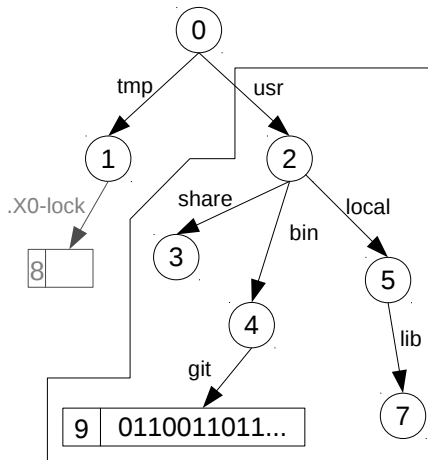
# Lock file specification

contains the lock file protocol



# Lock file specification

with other protocols



## Lock File specification

then, we can use the lock file specification:

$$\text{lock}(\text{path}) \sqsubseteq \langle \text{Lock}(\text{path}, v), \text{Lock}(\text{path}, 1) * v = 0 \rangle$$

$$\text{unlock}(\text{path}) \sqsubseteq \langle \text{Lock}(\text{path}, 1), \text{Lock}(\text{path}, 0) \rangle$$

## Lock File specification

$$\forall \text{path}, \boxed{P}. \text{islock}(\text{path}, \boxed{P}) \Rightarrow$$

$$\text{lock}(\text{path}) \sqsubseteq \langle \text{Lock}(\text{path}, v), \text{Lock}(\text{path}, 1) * v = 0 \rangle$$

$$\text{unlock}(\text{a}) \sqsubseteq \langle \text{Lock}(\text{path}, 1), \text{Lock}(\text{path}, 0) \rangle$$

$$\text{islock}(\text{path}, \boxed{P}) \triangleq \exists R. \boxed{P \iff \text{LF}(\text{path}) * R}$$

# Conclusions

- ▶ Introduced multi-atomic specifications
- ▶ Ordered sequences of atomic actions
- ▶ Unordered parallel atomic actions
- ▶ Formalised a fragment of POSIX file system operations
- ▶ Client reasoning
- ▶ Ownership in file systems is *cooperative*

# Future Work

- ▶ Link with operational models & testing
- ▶ Verify implementations
- ▶ Explore connection with refinement & Hoare's algebraic laws
- ▶ Mechanisation