

USC Viterbi

School of Engineering

Information Sciences Institute

Modeling cognitive workload in open-source communities via simulation

Alexey Tregubov, Jeremy Abramson, Christophe Hauser, Alefiya Hussain, and
Jim Blythe

Motivation

- Protecting open-source software (OSS)
- Understanding underlying socio-technical process
- Improve work experience of people engaged in the OSS community (e.g. Linux Kernel)



LKML

- Linux Kernel Mailing List (LKML) - OSS development artifacts: *code* & *communication*
- Developers, maintainers, reviewers discuss updates (patches) of the Linux Kernel.
- LKML messages capture a patch life cycle.

Last 100 messages	Today's messages	Yesterday's messages	Hottest Messages
Latest kernels			
mainline 6.4-rc2 patch			
stable 6.3.2 patch log			
stable 6.2.15 patch log			
longterm 6.1.28 patch log			
longterm 5.15.111 patch log			
longterm 5.10.179 patch log			
longterm 5.4.242 patch log			
longterm 4.19.282 patch log			
longterm 4.14.314 patch log			
Latest messages		Hottest messages	
Johan Hovold	Re: [PATCH v8 3/9] usb: dwc3: core: Access XHCI ad...	Linus Torvalds	Linux 6.4-rc2
Borislav Petkov	Re: [PATCH] EDAC: Expose node link in sysfs if CON...	Davidlohr Bueso	Re: Futexes & Folios
Niklas Schnelle	[PATCH v4 27/41] power: add HAS_IOPORT dependence...	Linus Torvalds	Re: [Regression w/ patch] Media commit causes user...
		Thomas Gleixner	Re: Futexes & Folios

LKML.ORG
Messages in this thread

- First message in thread
- Linus Torvalds
 - Jens Axboe
 - Linus Torvalds
 - Linus Torvalds
 - Linus Torvalds
- Mike Christie
 - Linus Torvalds
 - Mike Christie
 - Christian Brauner

[\[lkml\]](#) [\[2023\]](#) [\[May\]](#) [\[15\]](#) [\[last100\]](#) [RSS](#)
Views: [\[wrap\]](#) [\[headers\]](#) [\[forward\]](#)

From Linus Torvalds <>
Date Mon, 15 May 2023 08:56:39 -0700
Subject Re: [PATCH v11 8/8] vhost: use vhost_tasks for worker threads

On Mon, May 15, 2023 at 8:52 AM Jens Axboe <axboe@kernel.dk> wrote:

>
> Only potential downside is that it does make file references more
> expensive for other syscalls, since you now have a shared file table.
> But probably not something to worry about here?

Would the vhost user worker user processes ever be otherwise single-threaded?

I'd *assume* that a vhost user is already doing its own threads. But maybe that's a completely bogus assumption. I don't actually use any of this, so...

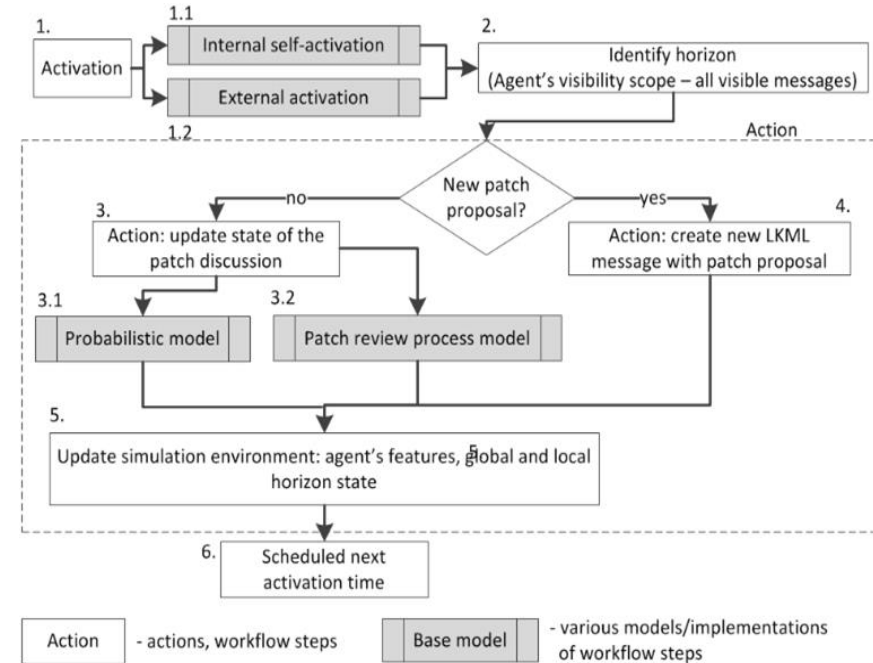
Because you are obviously 100% right that if you're otherwise single-threaded, then a CLONE_FILES kernel helper thread will cause the extra cost for file descriptor lookup/free due to all the race prevention.

Goals

- Develop an agent-based model of OSS development
- Explore benefits of explicit process models as part of multi-agent systems (LKML patch review process)
- Compare predictions of buggy/reverted LKML patches by multi-agent simulation and by ML models (ABS vs. ML predictions)
- Explore relationships the number of new patch proposals and the percent of patches that are later reverted

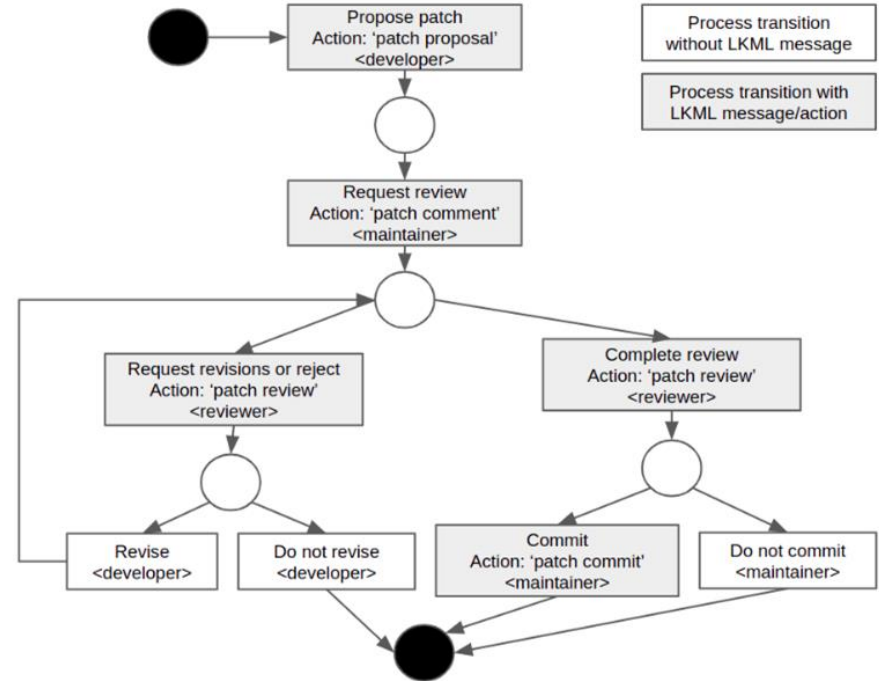
Agent-based Simulation Model of LKML

- DASH framework - a Python-based DES framework for modeling cognitive agents
- Patch discussions (message threads about patches) modeled explicitly, their status is updated by agents.
- Two simulation models:
 - probabilistic model
 - patch review process model



Patch review process

- Patch review process is modeled explicitly in our DASH model.
- State transitions are performed by agents/users with appropriate roles.



Simulation of Patch Reversal

In simulation models, the probability of patch reversal is dependent on likelihood of the developer to propose a patch that would be reverted (based on training data statistics) and of the reviewer to miss bugs in the review process.

We estimate the probability of patch being reversed as follows:

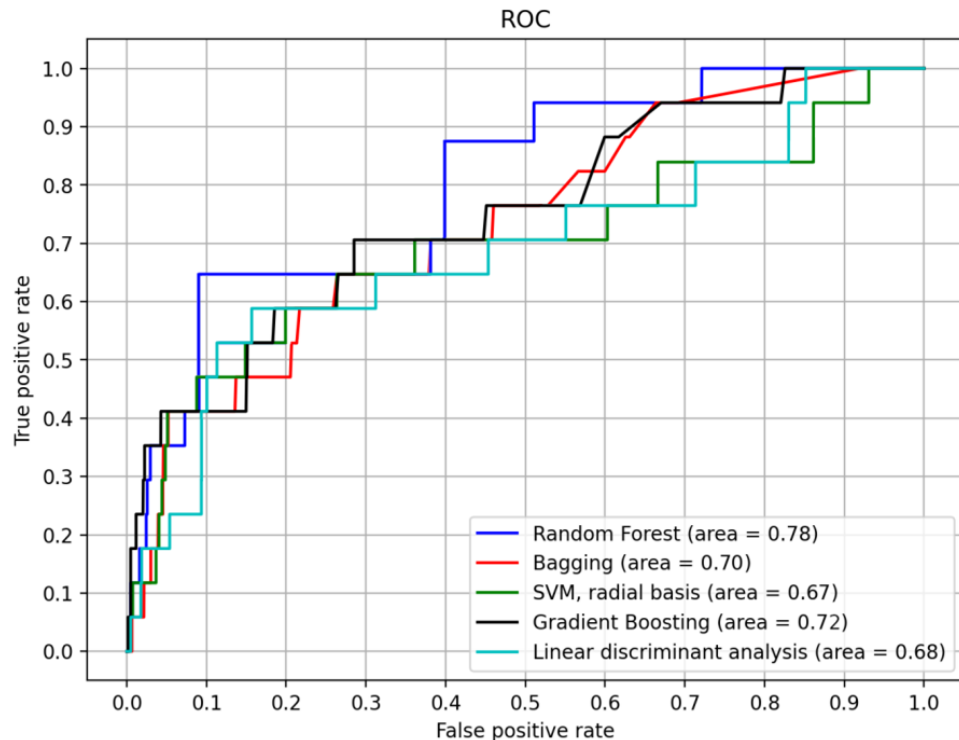
$$P = M_{load} \times P_{dev}$$

$$P_{dev} = \frac{n_patches_reverted}{n_patches_proposed}$$

M_{load} here is a reviewer workload/multitasking factor, which is a heuristic function that associates length of the work queue with increase of probability to miss a bug in the review process.

ML models for patch reversal prediction

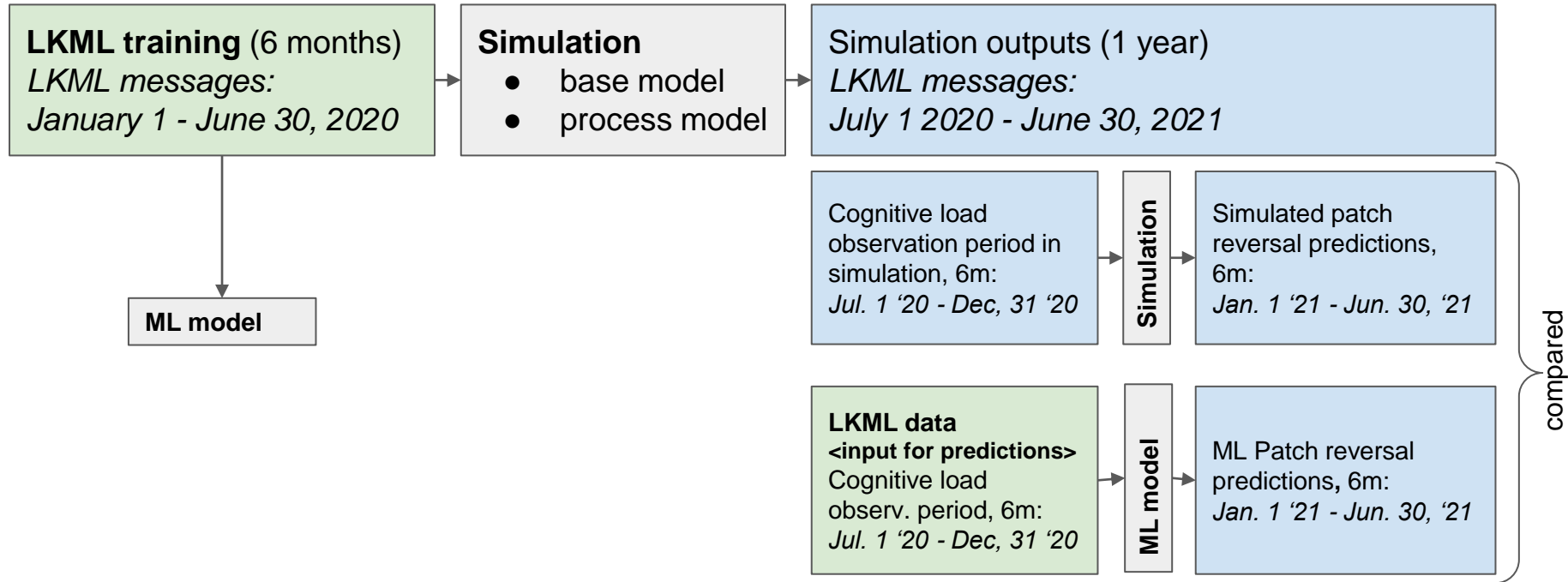
- Patch reversal prediction is a binary classification problem for a given feature vector of a LKML patch discussion
- Two groups of features used: code features and social features
- Several classifiers were compared, *Random Forest* classifier has highest AUC (0.78)



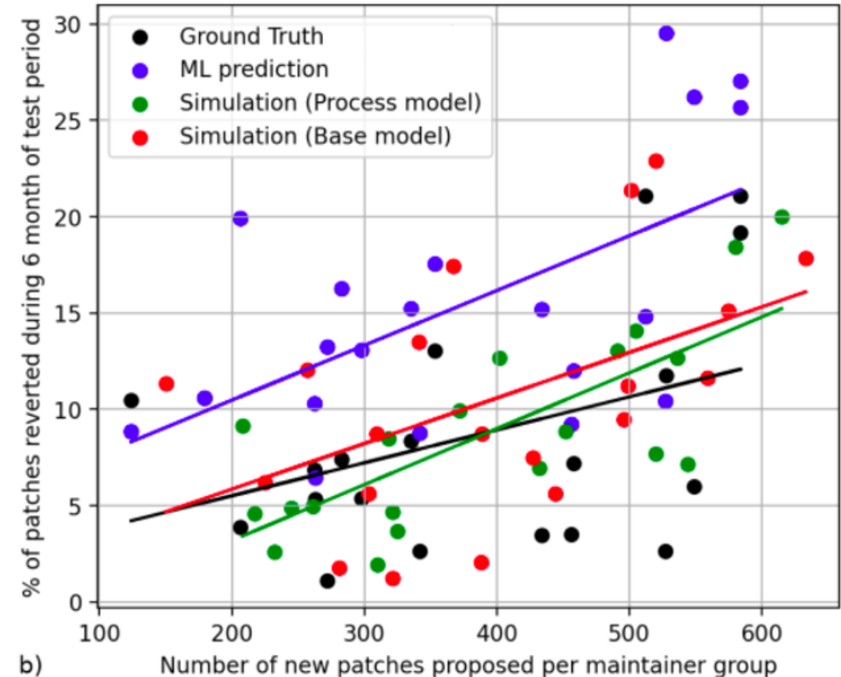
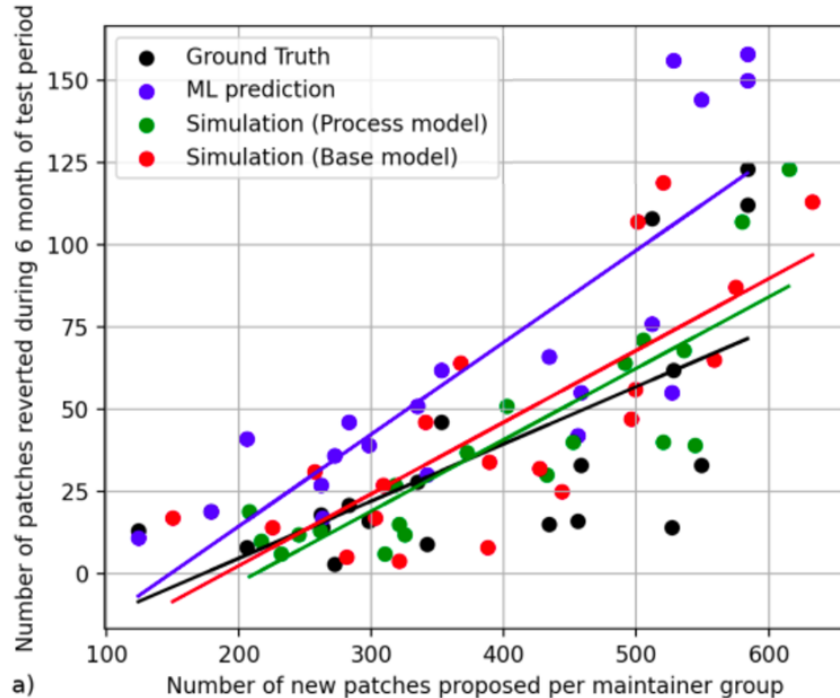
Experiment Setup

- Training: 212k LKML messages, 5K authors, January 1 - June 30 of 2020.
- Simulation: July 1 2020 - June 30 2021.
- The two simulation models: base model and process model
- Predictions on patch reversal were compared with ML model predictions.

Experiment Setup (cont-d)



Results: relation between number of new patches and patches reverted

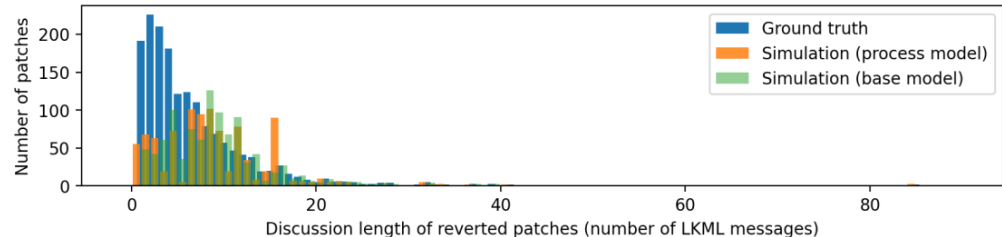


Results: goodness of fit to ground truth

- K-S test shows goodness of fit simulation models' predictions to the ground truth.

	D statistic	p-value
ML model	0.02107	0.7554
Simulation (Process model)	0.01618	0.9524
Simulation (Base model)	0.02157	0.7298

- Patch discussions produced by the simulation have different distributions of discussion lengths.



Contributions Summary

- Agent-based model of OSS development, including code review process
- Process model that captures repeated, structured multi-agent activities
- Multi-agent simulation improves the prediction of buggy (reverted) code in Linux Kernel over to statistical ML models
- Demonstrated that a simple model of cognitive workload shows the positive association between *the number of new patch proposals* and *the percent of patches that are later reverted**

*This association is captured by both the simulation model and ML models as well as observed in the ground truth.

Discussion, Q&A