

Tackling Deep Software Variability Together

Mathieu Acher @acherm

Special thanks to Luc Lesoil, Jean-Marc Jézéquel, Juliana Alves Pereira, and Paul Temple







Special thanks to Luc Lesoil, Jean-Marc Jézéquel, Juliana Alves Pereira, and Paul Temple



SOFTWARE VARIANTS ARE EATING THE WORLD AND SCIENCE







Learning variability can be easy...

State-of-the-art (SOTA) results: with the right algorithm and/or sampling strategy, fairly accurate and interpretable models for many systems and qualitative/quantitative properties can be derived

Still, why learning-based, SOTA techniques are not fully adopted and integrated to software-intensive projects?

- Academic/industry gap?
- Too costly?
- Hard to automatically observe/quantify variants?

References	Name	Domain	Non-Functional Properties
[61, 87]	Wget	Data transfer	memory footprint, code complexity
106]	Actian Vector	Database system	runtime
[41]	Apache Cassandra	Database system	latency
27, 29, 42, 48, 61,	Berkelev DB	Database system	L/O time, memory footprint, perfor-
2, 66, 80, 84-89, 98,			mance, response time, code com-
18, 119]			plexity, maintainability, binary size
[89]	FAME-DBMS	Database system	maintainability, binary size, perfor-
			mance
93 106 117 1201	MuSOI	Database system	defects throughout latency
106]	Rostaros	Database system	throughput latency
22 001	Demondar	Database system	monomy factorint performance
00-00 27 20 20 47 41	Crevayler	Database system	memory tootprint, performance
27, 29, 39, 47, 61,	SQLIE	Database system	memory tootprint, performance, re-
12, 00, 00-07, 07, 00,			sponse time, code complexity, run-
a, 105]	- 1 1	Property Research Contractory	time
15]	StockOnline	Database system	response time
10]	Katka	Distributed systems	throughput
25]	DNN	DNNs algorithms	accuracy of predictions
3]	Curriculum vitae	Document	number of pages
3]	Paper	Document	number of pages
[02	RUBiS	E-commerce application	response time
91]	EMAIL	E-mail client	time
47]	MBED TLS	Encryption library	response time
110]	SAP ERP	Enterprise Application	response time
63]	noc-CM-log	FPGA	CPU power consumption, runtime
63]	sort-256	FPGA	area, throughput
23]	E-Health System	Health	response time
29, 42, 84, 98]	HIPACC	Image processing	response time
17]	Disparity SPL	Image processing	energy consumption
86-88]	PKIab	Instant messenger	memory footprint, performance
36]	IBM ILOG CPLEX	Integer solver	runtime
110]	SPECiibb2005	Java Server	response time, throughput
101]	WEKA	Learning algorithm	accuracy of predictions
18]	SVD	Linear algebra	execution time and accuracy
42]	Trimach	Mash solver	iterations, menouse time
01	MRENCH	Miero henohmark	time
71) 72 1171	ACE. TAO mutant	Micro benchinark	dafaata
75, 116j ez eel	ACE+TAO system	Mature software	defects
00-001	Sensornetwork	Network simulator	memory tootprint, performance
109]	Sumonstrator	Network simulator	tatency
121]	NOC	Network-based system	energy and runtime
18]	Helmholtz 3D	Numerical analysis	execution time and accuracy
18]	Poisson 2D	Numerical analysis	execution time and accuracy
86-88, 92]	Linux kernel	Operating system	memory footprint, performance
40]	DNN	Optimization algorithm	response time
55]	Art system	Paint	user likeability
26]	Multigrid system	Equations solving	time of each interaction
[41]	CoBot System	Robotic system	CPU usage
42, 84, 98]	JavaGC	Runtime environment	response time
96]	Robot	Runtime environment	energy consumption and execution
			time
01-03, 00, 80, 1	54.		time. Feak Signai to Noise Ratio, re-
85, 87, 98, 105, 1	18.		sponse time, code complexity, video
119]			quality, performance
[98]	OpenCV	Video tracking	performance
[77]	C60 and MX300	Virtual environment	defects
[6]	Amazon EC2	Web coud service	performance
[27, 29, 48, 61, 6	52, Apache	Web server	response rate, response time, work-
66, 80, 84, 85, 9	98.		load, defects, throughput
117, 120]			6
[67]	Stock Brokerage	Web system	throughput, response time
[93]	vsftpd	FTP daemon	defects
[93]	ngIRCd	IRC daemon	defects

Learning variability can be easy...



... but we're not looking <u>deeply</u>.

Prediction models and configuration knowledge may not <u>generalize</u> (ie misleading/inaccurate/pointless for developers and users) if the software evolves, is compiled differently, a different hardware is used, input data fed differs, etc.



software application variability











deep software variability

hardware variability

15,000+ options

thousands of compiler flags and compile-time options

dozens of preferences

100+ command-line parameters

1000+ feature toggles

8



DEEP VARIABILITY

Sometimes, variability is consistent/stable across layers and knowledge transfer is immediate.

But there are also interactions among variability layers and <mark>variability knowledge</mark> may not generalize

Perf. 🛛

Bug

Perf. /

Poppler

264

node

🖹 X2

REAL WORLD EXAMPLE (X264)



REAL WORLD EXAMPLE (X264)



REAL WORLD EXAMPLE (X264)



DEEP VARIABILITY



The "best"/default software variant might be a bad one.

Influential software options and their interactions vary.

Performance prediction models and <mark>variability</mark> knowledge may not generalize



Bug

PERF. /

L. Lesoil, M. Acher, A. Blouin and J.-M. Jézéquel, "Deep Software Variability: Towards<mark>ump</mark> Handling Cross-Layer Configuration" in VaMoS 2021

Transferring Performance Prediction Models Across Different Hardware Platforms Valov et al. ICPE 2017

what about

variability of

nput data

Table 3: Sun ber of featur systems wer sured configu	nmary of es; NM – e measu urations	meas Nun red;	sured aber o NMC	systems; f machine – Numb	N _f – Num- es on which er of mea-
U	System	N_f	NM	NMC	

System	2 V J	1 8 1 9 1	TATALO .
XZ	7	7	154
x264	7	11	165
SQLite	5	10	32

Table 2: Summary of hardware platforms on which configurable software systems were measured; MID – Machine ID in DataMill cluster: NC – Number of CPUs; IS - Instruction set; CCR - CPU clock rate (MHz): RAM – RAM memory size (MB)

	Syster	ms		577001	Machine	es	
XZ	x264	SQLite	MID	NC	IS	CCR	RAM
1			73	2	i686	1733	1771
1	~	1	75	2	i686	3200	977
1			77	2	i686	2992	2024
1			78	1	i686	1495	755
~			79	4	x86_64	3291	7961
~			80	8	x86_64	3401	7907
~	~		81	16	x86_64	2411	32193
	~		87	1	i686	1595	249
	~		88	1	i686	1700	978
		1	90	2	i686	3200	977
	~		91	1	i686	2400	1009
	\checkmark	1	97	2	i686	2992	873
	~	1	98	2	i686	2992	873
		~	99	2	i686	2793	880
	~		103	2	i686	3200	881
	~		104	1	i686	1800	502
	~	~	105	2	i686	3200	881
	~		106	2	i686	3192	494
		✓	125	4	x86_64	3301	7960
	~		128	2	i686	2993	2024
		~	130	2	i686	3198	880
		~	146	2	i686	2998	872
		~	157	36	x86_64	2301	15954

"Linear model provides a good approximation of transformation between performance distributions of a system deployed in different hardware environments"

..... CPU HARDWARE TIT # CORES GPU OPERATING 1Pg SYSTEM OPTION ERSION DISTRIB ບາບ 264 ••• SOFTWARE COMPIL VARIAN VERSION compile-time options? INPUT DATA SIZE ENGTH

Transfer Learning for Software Performance Analysis: An Exploratory Analysis Jamshidi et al. ASE 2017

SPEAR (SAT Solver) Analysis time	X264 (video enc Encoding time	ode	r)	SQLite (DB engine) Query time	
14 options 16,384 configurations SAT problems 3 hardware 2 versions	16 options 4,000 configurations Video quality/size 2 hardware			14 options 1,000 configurations DB Queries 2 hardware	
2 Versions	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	SM S M L L L VL	0.97 0.96 0.65 0.67 0.05 0.06 0.08 0.09		

Insight. For non-severe hardware changes, we can linearly transfer performance models across environments.

Insight. The strength of the influence of configuration options is typically preserved across environments.

Insight. A large percentage of configurations are typically invalid in both source and target environments.

SaC (Compiler) Execution time

50 options 71,267 configurations 10 Demo programs



"Analyzing the Impact of Workloads on Modeling the Performance of Configurable Software Systems" Stefan Mühlbauer et al. ICSE 2023

Summary (RQ_1): Varying the workload causes a substantial amount of variation among performance distributions. Across workloads, we observed *mostly linear* (for six of the nine subject systems), but to a large extent, also *non-monotonous* relationships (for three of the nine subject systems).

Summary (RQ_2): Workloads can affect performance influences of configuration options in various ways (e.g., conditioning influence, introducing variance, having outliers). We can correlate some variation of performance influences with workload characteristics, yet identifying relevant workload characteristics is highly domain-specific and cannot be considered trivial.



(c) Distribution of configuration options' performance influences under different workloads for H2. Table IV: Common top five influential configuration options among pairs of workloads.

System	Workload 1	Workload 2	#Common
JUMP3R	helix.wav	sweep.wav	2
KANZI	vmlinux	fannie_mae_500k	1
DCONVERT	jpeg-small	svg-large	2
н2	tpcc-2	tpcc-8	3
BATIK	village	cubus	4
XZ	deepfield	silesia	4
LRZIP	artificl	uiq-32-bin	3
x264	sd_crew_cif_short	sd_city_4cif_short	4
z3	QF_NRA_hong_9	QF_BV_bench_935	3



Let's go deep with input data!



Intuition: video encoder behavior (and thus runtime configurations) hugely depends on the input video (different compression ratio, encoding size/type etc.)

- Is the best software configuration still the best?
- Are influential options always influential?
- Does the configuration knowledge generalize?

two performance models f_1 and f_2 $f_1 = \beta \times f_2 + \alpha$



YouTube User General Content dataset: **1397 videos** Measurements of **201 soft. configurations** (with <u>same</u> hardware, compiler, version, etc.): encoding time, bitrate, etc.



Do x264 software performances stay consistent across inputs?

- Encoding time: very strong correlations
 o low input sensitivity
- FPS: very strong correlations

 $\,\circ\,$ low input sensitivity

- CPU usage : moderate correlation, a few negative correlations
 medium input sensitivity
- Bitrate: medium-low correlation, many negative correlations
 - High input sensitivity
- Encoding size: medium-low correlation, many negative correlations
 High input sensitivity

two performance models f_1 and f_2 $f_1 = \beta \times f_2 + \alpha$

1397 videos x 201 software configurations





Are there some configuration options more sensitive to input videos? (P = bitrate)





L. Lesoil, M. Acher, A. Blouin and J.-M. Jézéquel "Input sensitivity on the performance of configurable systems an empirical study" JSS 2023

Practical impacts for users, developers, scientists, and self-adaptive systems

Threats to variability knowledge: predicting, tuning, or understanding configurable systems without being aware of inputs can be inaccurate and... pointless

eg effectiveness of sampling strategies (random, 2-wise, etc.) is input specific (see also Pereira et al. ICPE 2020)

Opportunities: for some performance properties (P) and subject systems, some stability is observed and performance remains consistent!

Stefan Mühlbauer, Florian Sattler, Christian Kaltenecker, Johannes Dorn, Sven Apel, Norbert Siegmund "Analyzing the Impact of Workloads on Modeling the Performance of Configurable Software Systems" ICSE 2023

Insight: Workload sensitivity challenges the robustness and $\frac{1}{30}$
generalizability of single-workload performance models, yet $\frac{100}{51}$
it is neglected in state-of-the-art approaches. Worse, robust 33
techniques using only rankings or relative importance of $\frac{180}{50}$
options are inapplicable for certain workload variations.

System	#M	Performance(s) P	Docker	Dataset
gcc	2400	size, ctime, exec	Link	Link
ImageMagick	100 000	size, time	Link	Link
lingeling	35 100	#confl.,#reduc.	Link	Link
nodeJS	96 950	#operations/s	Link	Link
poppler	23 680	size, time	Link	Link
SQLite	7500	15 query times q1-q15	Link	Link
x264	280 797	cpu, fps, kbs, size, time	Link	Link
xz	1440	size, time	Link	Link

L. Lesoil, M. Acher, A. Blouin and J.-M. Jézéquel "Input sensitivity on the performance of configurable systems an empirical study" JSS 2023



Is there an interplay between compile-time and runtime options?



L. Lesoil, M. Acher, X. Tërnava, A. Blouin and J.-M. Jézéquel "The Interplay of Compiletime and Run-time Options for Performance Prediction" in SPLC '21



This paper investigates how compile-time options can affect software performances and how compile-time options interact with run-time options.

Figure 1: Cross-layer variability of x264





variants and versions: The case of linux kernel size" Transactions on Software Engineering (TSE), 2021





What can we do? (#1 studies)

Empirical studies about deep software variability

- more subject systems
- more variability layers, including interactions
- more quantitative (e.g., performance) properties

with challenges for gathering measurements data:

- how to scale experiments? Variant space is huge!
- how to fix/isolate some layers? (eg hardware)
- how to measure in a reliable way?

Expected outcomes:

- significance of deep software variability in the wild
- identification of stable layers: sources of variability that should not affect the conclusion and that can be eliminated/forgotten
- identification/quantification of **sensitive** layers and interactions that matter
- variability knowledge



What can we do? (#2 cost)

Reducing the cost of exploring the variability spaces

- learning
 - many algorithms/techniques with tradeoffs interpretability/accuracy
 - transfer learning (instead of learning from scratch) Jamshidi et al. ASE'17, ASE'18, Martin et al. TSE'21
- sampling strategies
 - uniform random sampling? t-wise? distance-based? ...
 - sample of hardware? input data?
- incremental build of configurations Randrianaina et al. ICSE'22
- white-box approaches Velez et al. ICSE'21, Weber et al. ICSE'21



What can we do? (#3 modelling)

Modelling variability

- Abstractions are definitely needed to...
 - reason about logical constraints and interactions
 - integrate domain knowledge
 - synthesize domain knowledge
 - automate and guide the exploration of variants
 - scope and prioritize experiments
- Challenges:
 - Multiple systems, layers, concerns
 - Different kinds of variability: technical vs domain, accidental vs essential, implicit vs explicit... when to stop modelling?
 - \circ reverse engineering



Open, reproducible science for deep variability

deep.variability.io?

A collaborative resource and place:

- Evidence of Deep Software Variability
- Continuous survey of papers/articles about deep software variability
- Case studies and configuration knowledge associated to software projects
- Datasets (eg measurements of performance in certain conditions)
- Reproducible scripts (eg for building prediction models)
- Description and results of Innovative solutions (eg transfer learning)
 - content based on already published papers
 - beyond PDFs

Ongoing work ("live" book with jupyter)... Looking for contributions/ideas/suggestions!

no replicat e

RCS
 RCS
 RCS

11 READVE.cd

Bangakanguane
 Bangakanguane
 Bangakanguane
 Bangakanguane

T 155 fram, Senior Millerry



JSS_Input_Sensitivity.p...



Learning variability can be easy...



G

Tackling Deep Software Variability <u>Together</u>

Mathieu Acher @acherm



BACKUP SLIDES

(BEYOND X264) EMPIRICAL AND FUNDAMENTAL QUESTION



HOW DOES DEEP SOFTWARE VARIABILITY MANIFEST IN THE WILD?

SOFTWARE SCIENTISTS SHOULD OBSERVE THE JUNGLE/ GALAXY!

DEEP SOFTWARE VARIABILITY

4 CHALLENGES AND OPPORTUNITIES

Luc Lesoil, Mathieu Acher, Arnaud Blouin, Jean-Marc Jézéquel: Deep Software Variability: Towards Handling Cross-Layer Configuration. VaMoS 2021: 10:1-10:8

IDENTIFY THE INFLUENTIAL LAYERS

1



TEST & BENCHMARK ENVIRONMENTS



PROBLEM

COMBINATORIAL EXPLOSION AND COST









<u>Challenge</u> Build a representative, cheap set of environments

Software





0.152.2854



0.155.2917









TRANSFER PERFORMANCES ACROSS ENVIRONMENTS



REDUCE COST OF MEASURE

(3)



CHALLENGES FOR DEEP SOFTWARE VARIABILITY

IDENTIFY THE INFLUENTIAL LAYERS

TEST & BENCHMARK ENVIRONMENTS

TRANSFER PERFORMANCES ACROSS ENVIRONMENTS

CROSS-LAYER TUNING

x264 video encoder (compilation/build)

x264 --bframes 1 --ref 3 --cabac DiverSE-meeting.mp4 -o meeting13.webm

mathieuacher

localhost.localdomain

./configure --help

--disable-thread --disable-win32thread --disable-interlaced --bit-depth=BIT DEPTH --chroma-format=FORMAT disable multithreaded encoding disable win32threads (windows only) disable interlaced encoding support set output bit depth (8, 10, all) [all] output chroma format (400, 420, 422, 444, all) [all]

disable platform-specific assembly optimizations

Advanced options:

--disable-asm --enable-lto --enable-debug

--enable-gprof

--enable-strip

--enable-pic

Cross-compilation:

- --host=HOST --cross-prefix=PREFIX
- --svsroot=SYSROOT

External library support:

- --disable-avs --disable-swscale --disable-lavf --disable-ffms
- --disable-gpac

--disable-lsmash

enable link-time optimization add -g add -pg add -s build position-independent code

build programs to run on HOST use PREFIX for compilation tools root of cross-build tree

disable avisynth support disable swscale support disable libavformat support disable ffmpegsource support disable gpac support disable lsmash support





Key results (for x264)

L. Lesoil, M. Acher, X. Tërnava, A. Blouin and J.-M. Jézéquel "The Interplay of Compiletime and Run-time Options for Performance Prediction" in SPLC '21

Worth tuning software at compile-time: gain about 10 % of execution time with the tuning of compile-time options (compared to the default compile-time configuration). The improvements can be larger for some inputs and some runtime configurations.

Stability of variability knowledge: For all the execution time distributions of x264 and all the input videos, the worst correlation is greater than 0.97. If the compile-time options change the scale of the distribution, they do not change the rankings of run-time configurations (i.e., they do not truly interact with the run-time options).

<u>Reuse of configuration knowledge</u>: $f_1 = \beta \times f_2 + \alpha$

- Linear transformation among distributions
- Users can also trust the documentation of run-time options,

consistent whatever the compile-time configuration is.





Are there some configuration options more sensitive to input videos? (bitrate)





H. Martin, M. Acher, J. A. Pereira, L. Lesoil, J. Jézéquel and D. E. Khelladi, "Transfer learning across variants and versions: The case of linux kernel size" Transactions on Software Engineering (TSE), 2021

- Linux as a subject software system (not as an OS interacting with other layers)
- Targeted non-functional, quantitative property: <u>binary size</u>
 - interest for maintainers/users of the Linux kernel (embedded systems, cloud, etc.)
 - challenging to predict (cross-cutting options, interplay with compilers/build systems, etc/.)
- Dataset: version 4.13.3 (september 2017), x86_64 arch, measurements of 95K+ random configurations
 - paranoiac about deep variability since 2017, Docker to control the build environment and scale
 - diversity of binary sizes: from 7Mb to 1.9Gb
 - 6% MAPE errors: quite good, though costly...





H. Martin, M. Acher, J. A. Pereira, L. Lesoil, J. Jézéquel and D. E. Khelladi, "Transfer learning across variants and versions: The case of linux kernel size" Transactions on Software Engineering (TSE), 2021

Transfer learning to the rescue

- Mission Impossible: Saving variability knowledge and prediction model 4.13 (15K hours of computation)
- Heterogeneous transfer learning: the feature space is different



• TEAMS: transfer evolution-aware model shifting





Joelle Pineau "Building Reproducible, Reusable, and Robust Machine Learning Software" ICSE'19 keynote "[...] results can be brittle to even minor perturbations in the domain or experimental procedure"

than the average.

Deep Reinforcement Learning that Matters

Peter Henderson^{1*}, Riashat Islam^{1,2*}, Philip Bachman² Joelle Pineau¹, Doina Precup¹, David Meger¹ What is the magnitude of the effect hyperparameter settings can have on baseline performance?

How does the choice of **network architecture** for the policy and value function approximation affect performance?

How can the **reward scale** affect results?

Can random seeds drastically alter performance?

How do the **environment properties** affect variability in reported RL algorithm performance?

Are commonly used baseline **implementations** comparable?





"Neuroimaging pipelines are known to generate different results depending on the computing platform where they are compiled and executed."

Reproducibility of neuroimaging

analyses across operating systems,

Glatard et al., Front. Neuroinform., 24

April 2015



	Cluster A	Cluster B
Applications	Freesurfer 5.3.0, build 1	Freesurfer 5.3.0, build 1 and 2
	FSL 5.0.6, build 1	FSL 5.0.6, build 1 and 2
	CIVET 1.1.12-UCSF, build 1	CIVET 1.1.12-UCSF, build 1
Interpreters	Python 2.4.3, bash 3.2.25,	Python 2.7.5, bash 4.2.47,
	Perl 5.8.8, tcsh 6.14.00	Perl 5.18.2, tcsh 6.18.01
glibc version	2.5	2.18
OS	CentOS 5.10	Fedora 20
Hardware	x86_64 CPUs (Intel Xeon)	x86_64 CPUs (Intel Xeon)

Statically building programs improves reproducibility across OSes, but small differences may still remain when dynamic libraries are loaded by static executables[...]. When static builds are not an option, software heterogeneity might be addressed using virtual machines. However, such solutions are only workarounds: differences may still arise between static executables built on different OSes, or between dynamic executables executed in different VMs.



DEEP QUESTIONS?