



EMORY
UNIVERSITY



Sandia
National
Laboratories



Beyond Guess and Check: Quantifying the Fidelity of Proxy Applications

Si Chen¹, Simon Garcia de Gonzalo², Omar Aaziz², Jeanine Cook², **Avani Wildani**^{1,3}

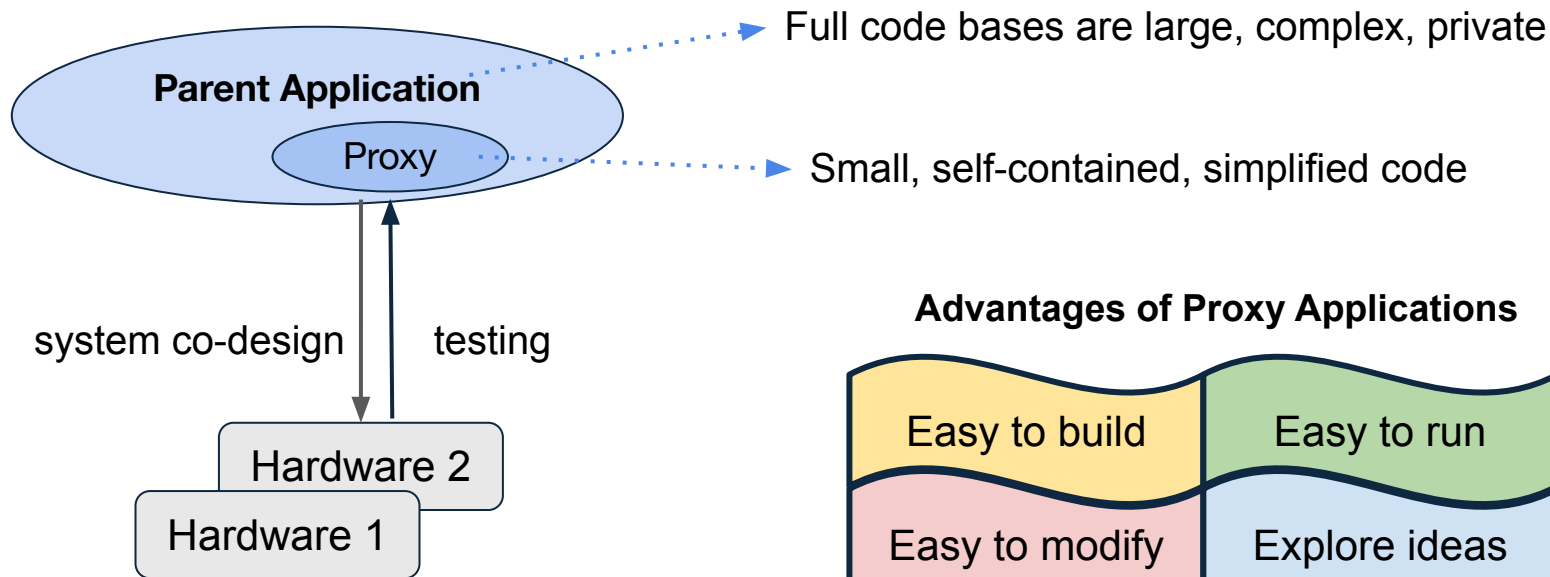
PMBS'25, held in conjunction with SC'25
November 2025

¹ Emory University

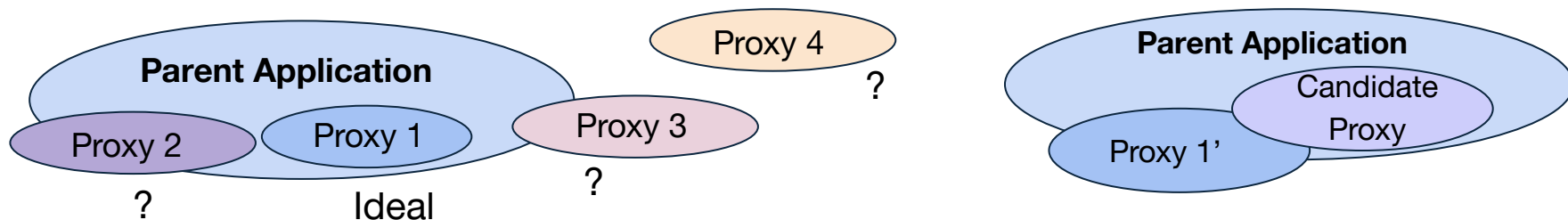
² Sandia National Laboratories

³ Cloudflare

Background: Proxy Applications for co-design



Motivation: Fidelity of Proxy Applications



Fidelity: the degree of accuracy with which a proxy application mimics the behavior and performance characteristics of its parent application.

Why quantify the fidelity? Poor fidelity lead to suboptimal architectural decisions

Motivation

Previous Work:

- Limited number of applications
- Limited features or metrics
- Qualitative analysis
- Lack of standard metrics for proxy

Three research questions:



How to quantify
similarity?



How to **select** a
representative
set of proxies?



How to
effectively
identify the
discrepancies?



Image: MoMa, Calder's Yellow Anteater

We introduce **Calder**, a tool kit that implements **similarity algorithms** to compare application behaviors, using **Laplacian scores** and a correlation filter to select the **most important** application behavior features.

Key Contributions

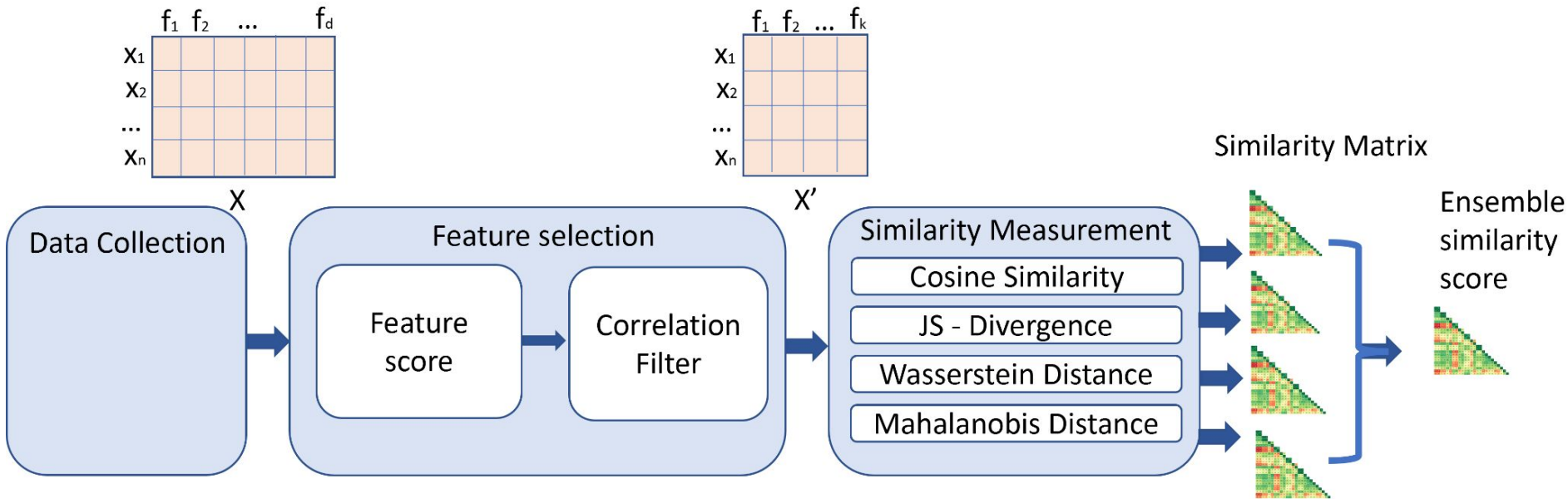
- **Quantitative characterization** of fidelity across a broad range of proxy and parent application pairs.
- Advanced **feature selection** techniques that reduce data size by up to 95%.
- Usage of kernel function similarity as ground truth to **define similarity** threshold.



Outline

- Background and Motivation
- Calder Architecture
- How to quantify the similarity?
- How to select a representative set of proxies?
- How to effectively identify the discrepancies?
- Conclusion

Calder Architecture



“Features” in this talk are hardware performance counters

Data Collection

- PAPI events collected by **LDMS** infrastructure.
- 600+ hardware events
- 15 subgroups according to event categories (e.g., Dispatch Pipeline, Instruction Cache)
- Input size: same problem or 50% memory usage
- MPI-only mode, 128 ranks on 4 nodes average, end time point accumulated value
- Normalized with CPU cycles

- **Hardware Platforms**
 - Intel Skylake
 - IBM Power9

To ensure robustness and reliability, each subgroup was run five times for each application, totaling over 3,000 runs.

Data Collection: Applications

Table 1: Proxy/Parent and Control Apps

Proxy	Parent	Other apps
ExaMiniMD	LAMMPS	AMG2013
miniQMC	QMCPACK	Castro
miniVite	Vite	Laghos
Nekbone	Nek5000	PENNANT
PICSARlite	PICSAR	SNAP
SW4lite	SW4	HPCG benchmark
SWFFT	HACC	HPCC benchmark
XSbench	OpenMC	

Data Collection: Applications

An application is represented by a long **vector** with multiple events.

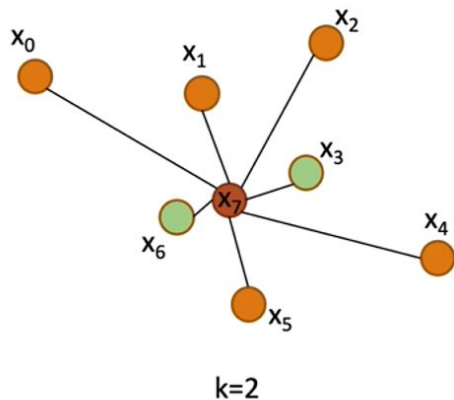
	BR_INST_RETIRED:ALL_BRANCHES	BR_INST_RETIRED:CONDITIONAL	BR_INST_RETIRED:FAR_BRANCH
→ ExaMiniMD	0.102941117	0.084736321	3.04E-06
→ LAMMPS	0.166040571	0.132466767	5.41E-06
→ sw4lite	0.066925525	0.042952159	1.14E-05
→ sw4	0.060973991	0.039039782	1.05E-05

...

...

Feature Selection: Laplacian Score

Laplacian score builds a **nearest neighbor graph** for application points and seeks those features that respect **local** graph structure to **preserve the similarity structure**.



Step 1: Derive the nearest neighbor graph

	X0	X1	X2	X3	X4	X5	X6	X7
X0								0
X1								0
X2								0
X3								$S_{3,7}$
X4								0
X5								0
X6								$S_{6,7}$
X7	0	0	0	$S_{7,3}$	0	0	$S_{7,6}$	0

Step 2: Calculate the Laplacian

$$L_r = \frac{\text{neighbor difference for feature } r}{\text{Variance of feature } r}$$

The **smaller** the Laplacian score is, the **more important** the feature is!



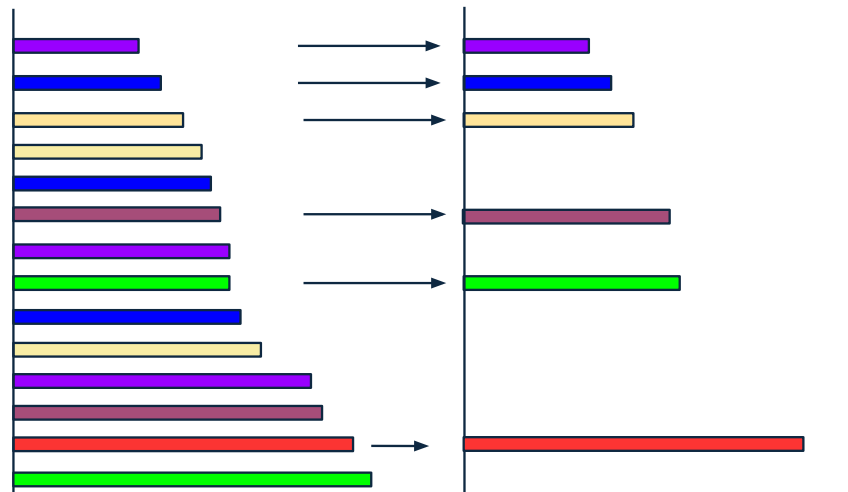
Feature Selection: Correlated Features

L ranks features by importance 😊...but does not consider the relationship between features. 😞

Pearson Correlation Coefficient (PCC)

measures the linear correlation between two variables:

$$\rho_{X,Y} = \frac{\mathbb{E}[(X - \mu_X)(Y - \mu_Y)]}{\sigma_X \sigma_Y}$$



Ranked features with
Laplacian scores

Ranked uncorrelated
features

Similarity and Distance

To evaluate the similarity between applications, we calculate the pairwise distance for each application pair using four representative similarity measurement methods.

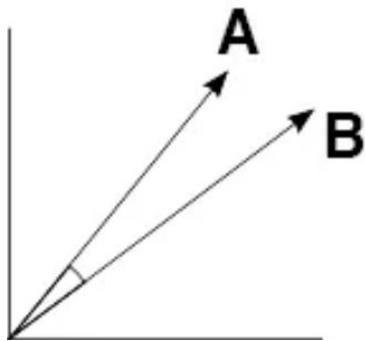
- Cosine similarity
- Jensen-Shannon (JS) divergence
- Wasserstein Distance (WD)
- Mahalanobis Distance (MaD)



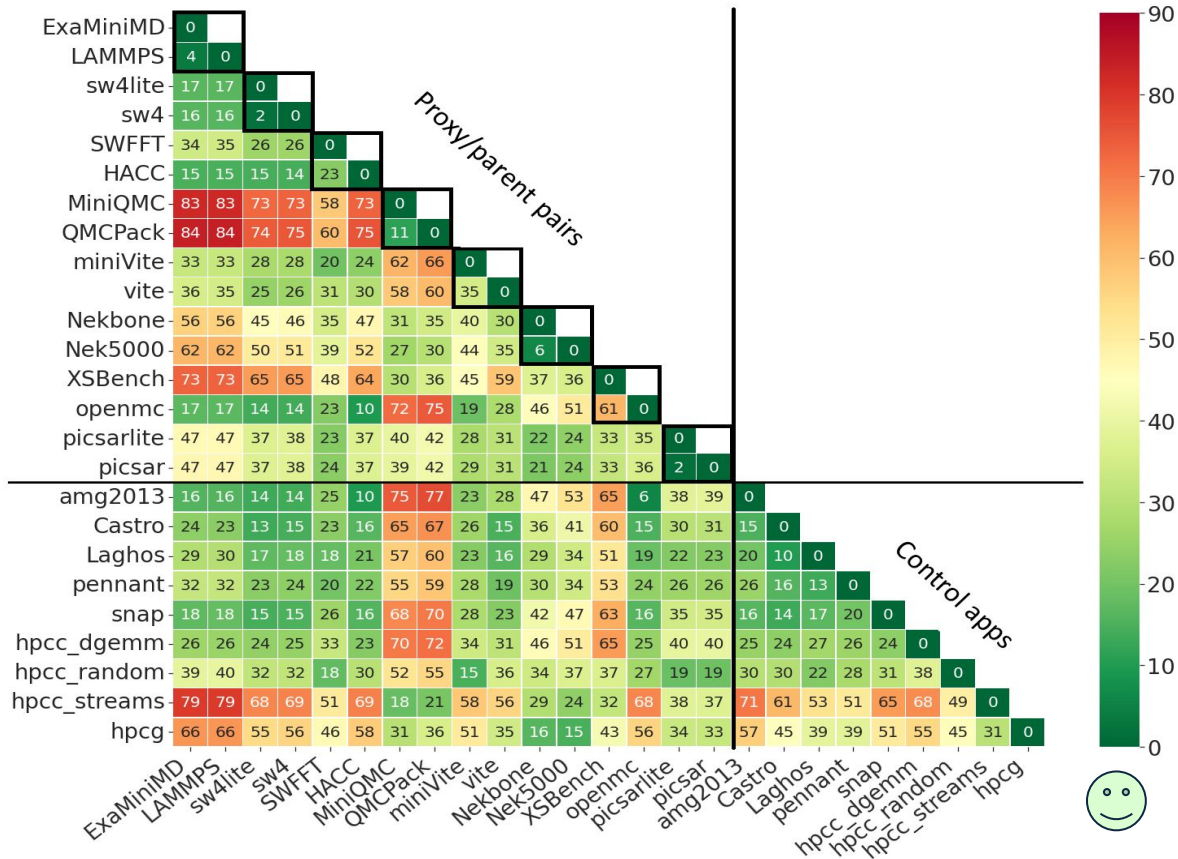
Cosine Similarity

Similarity score

0	
<30	0



Intel Skylake



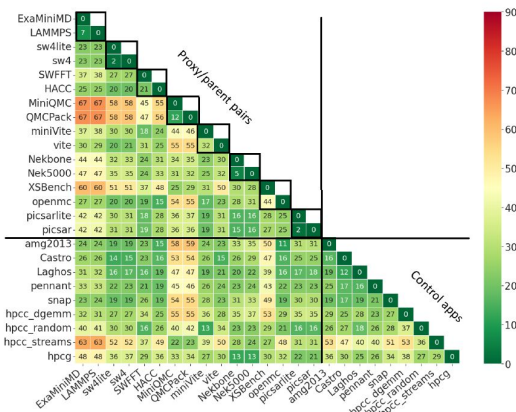


How to quantify similarity?

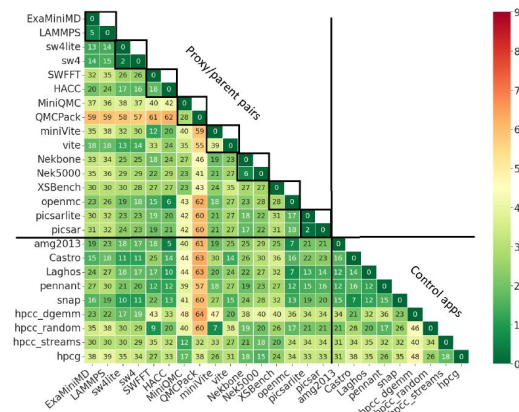
Similar proxy/parent pairs maintain consistency across similarity algorithms.

Dissimilarities are algorithm-dependent.

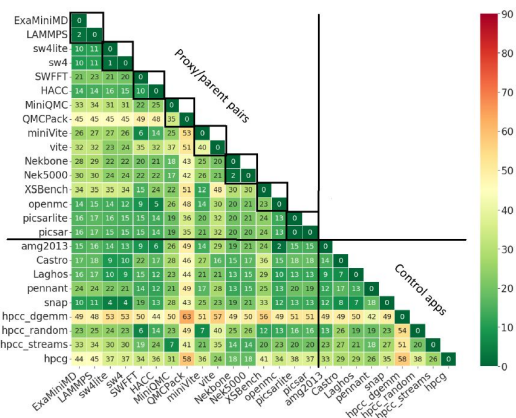
We select **cosine similarity** for validating fidelity.



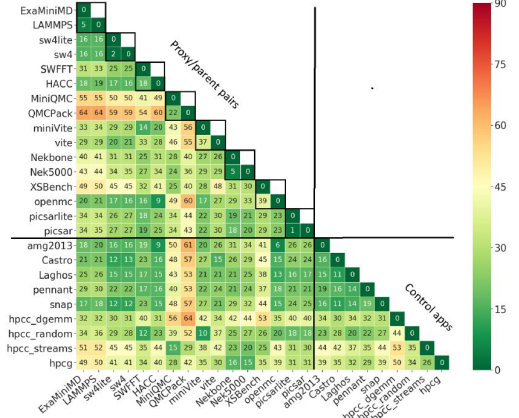
(a) JS divergence (100x)



(b) Wasserstein distance (1000x)



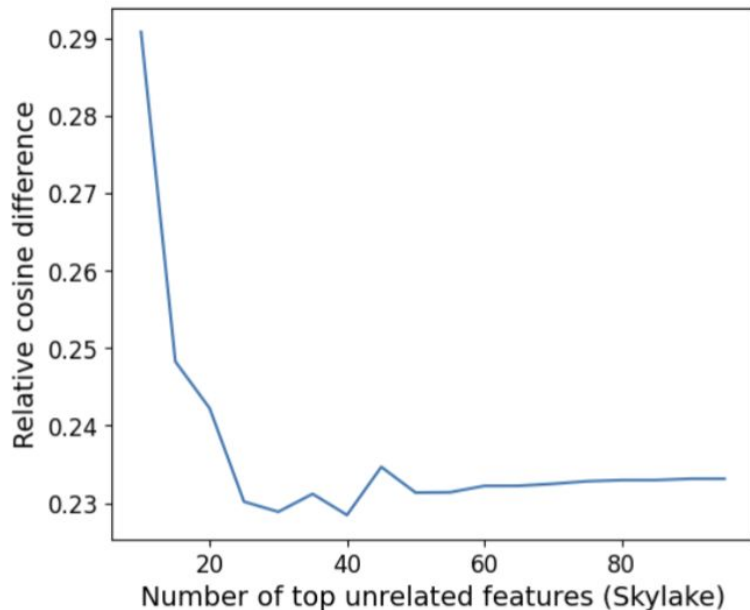
(c) Mahalanobis distance (10x)



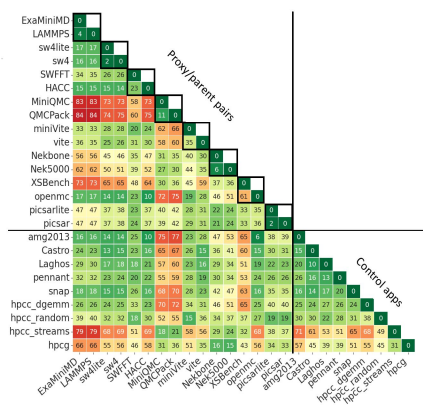
(d) Similarity for ensemble methods



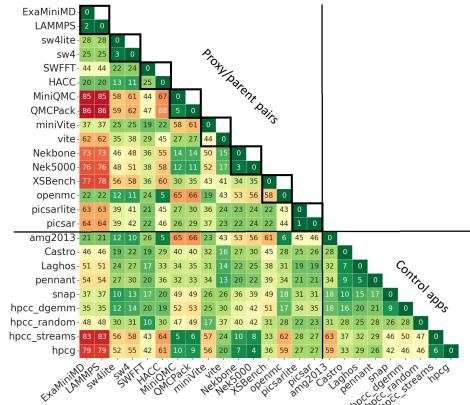
Feature sensitivity and Top Features



Cosine Similarity



All features



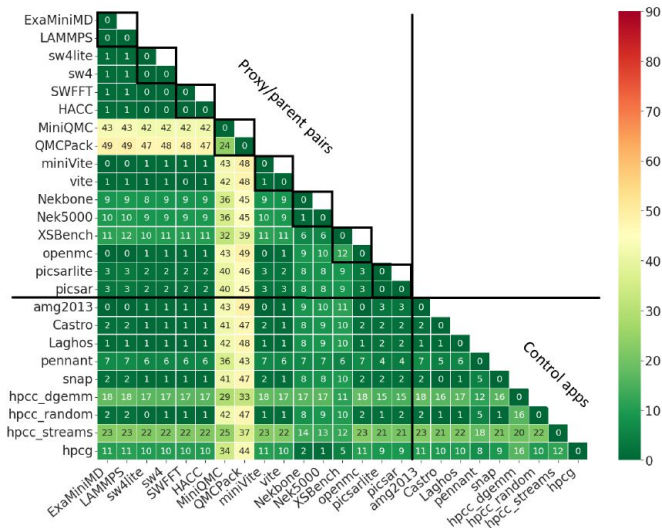
Top 25 unrelated features



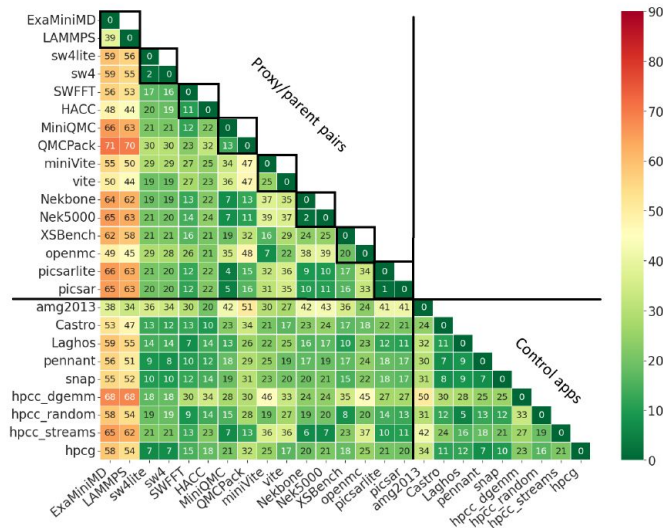
Subgroup Features Similarity

How to **select** a representative set of proxies?

Help proxy developers to observe the behavior differences to tackle code changes better and produce sounder proxies.



(a) L1 Cache



(b) Memory Pipeline



Feature Standard Deviation

The variability of hardware
performance counters (features)
during runtime

Substantial feature-level
deviations often reflect deeper
behavioral divergence.

How to effectively
identify the
discrepancies?

Table 4: Dissimilarity Feature Source for Proxy/Parents Pairs

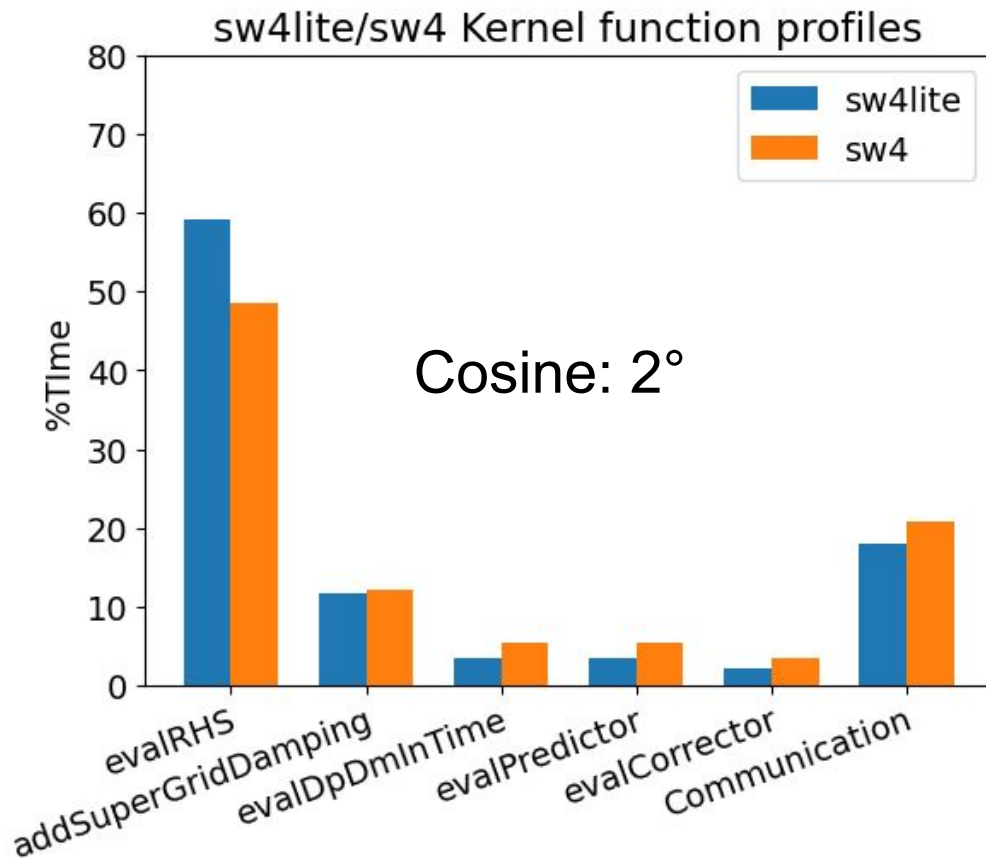
Proxy and Parent pairs	>2std	>3std	>4std	>5std
ExaMiniMD / LAMMPS	10	8	8	4
SW4lite / SW4	1	1	1	1
SWFFT / HACC	17	12	11	8
miniQMC / QMCPACK	13	10	8	6
miniVite / Vite	72	38	23	21
Nekbone / Nek5000	11	6	2	2
XSbench OpenMC	16	9	9	8
PICSARlite / PICSAR	1	1	1	0
Unique feature #s	99	64	49	38



Validating Similarity: Root Cause Analysis

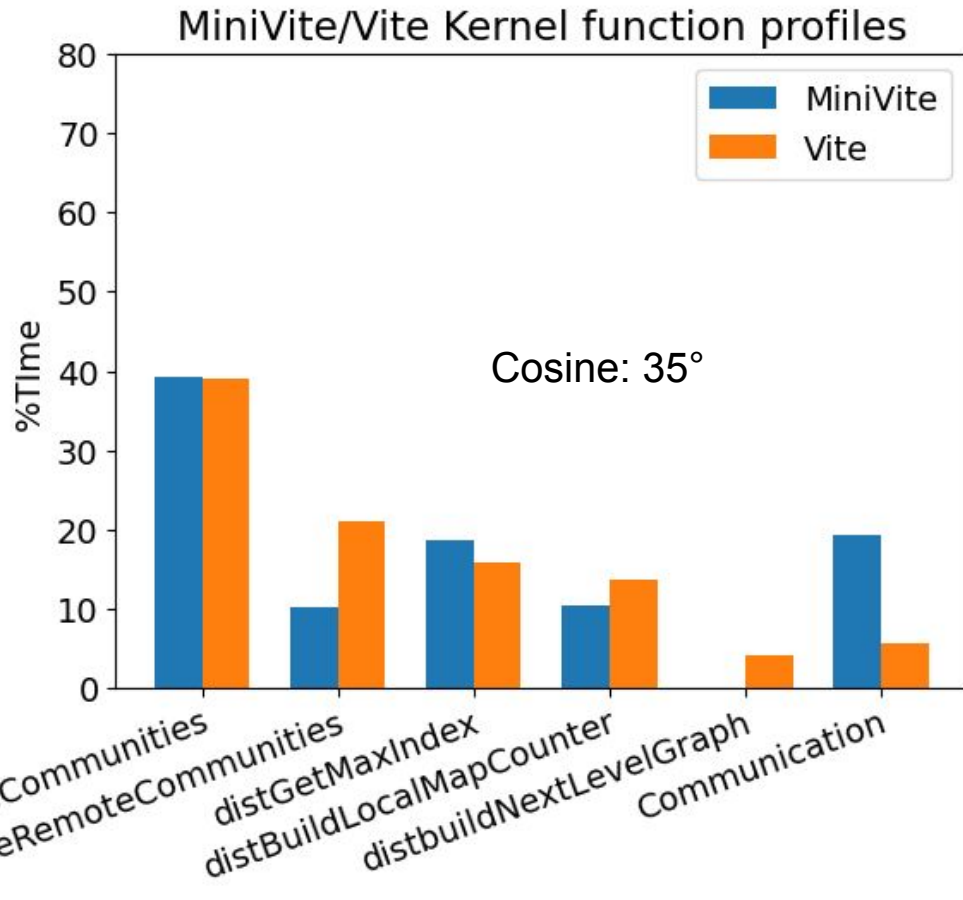
Ground truth for proxy-parent similarity:
code base implementations

Kernel function profiles (normalized
execution time) collected with gporf.



Less Rosy Root Cause Analysis

We recommend a general **similarity threshold** of **30°** for cosine similarity.



Conclusion

Calder **identifies** representative features, **minimizing** data collection overhead, and **highlights** key dissimilarities to **guide** future proxy co-design.



How to quantify
similarity?

Cosine Similarity



How to **select** a
representative set
of proxies?

**Subgroup
Features**

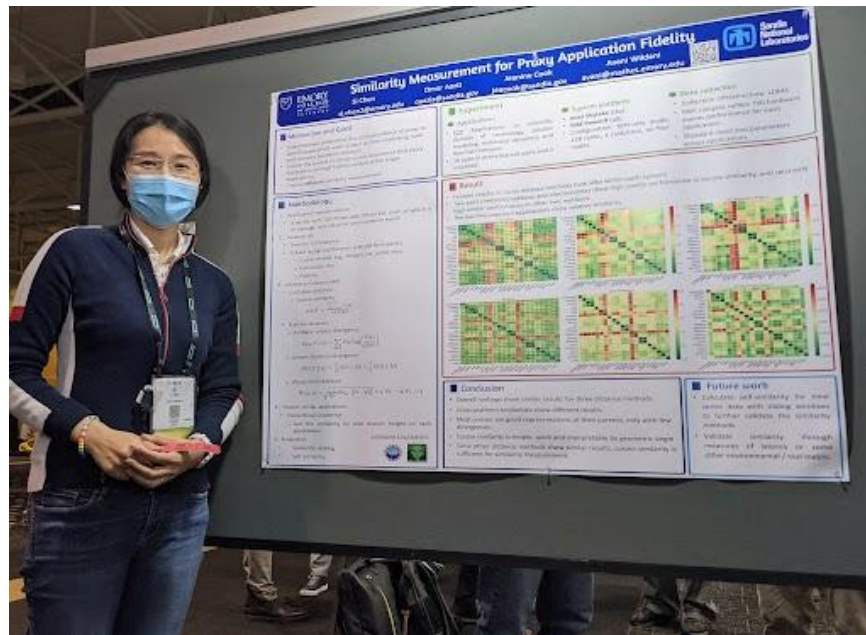
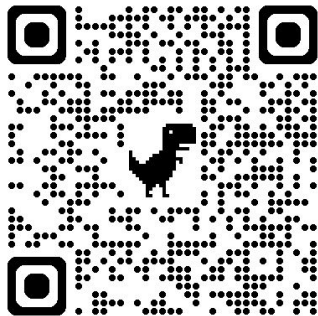


How to effectively
identify the
discrepancies?

Dissimilarity

Thank you!

GitHub



CV of Dr. Si Chen



meditates@gmail.com

