

A Multi-objective Adaptivity Methodology Vertically Integrating Algorithmic Parameters and Design Space Exploration

Luigi Nardi, PhD

Software Performance Optimisation group
@Argonne National Lab

January 16th 2017

In collaboration with:

B. Bodin, M Z. Zia, H. Wagstaff, D. Carroll, A. White, E. Vespa, S. Saeedi, G. S. Shenoy, M. K. Emani, J. Mawer, A. Nisbet, M. Luján, B. Franke, G. Riley, M. F. P. O'Boyle, A. J. Davison, P. H. J. Kelly and S. Furber



The University of Manchester

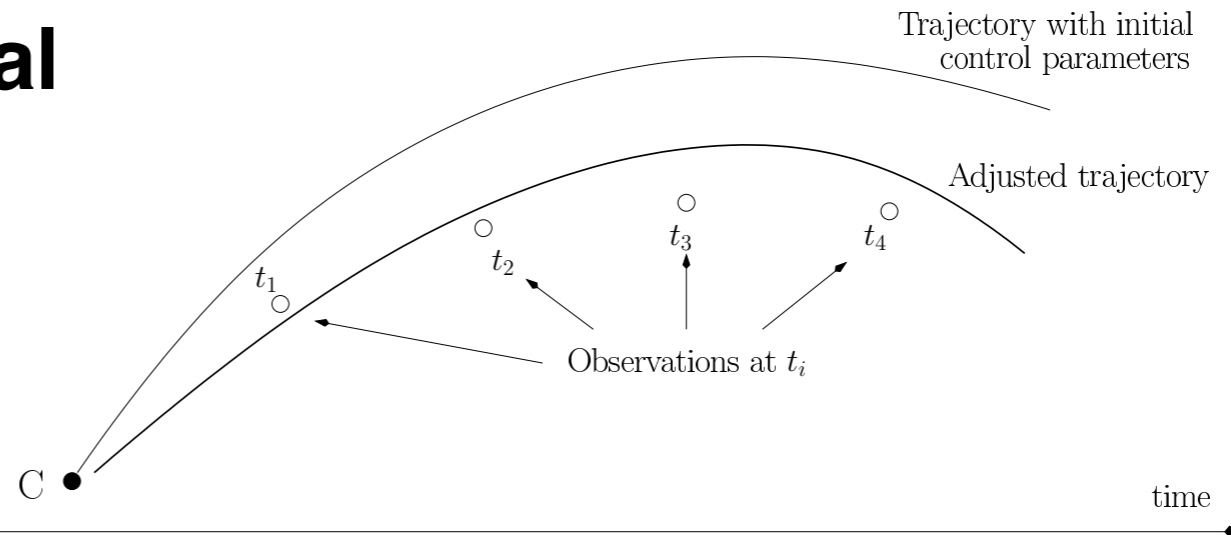


Imperial College
London



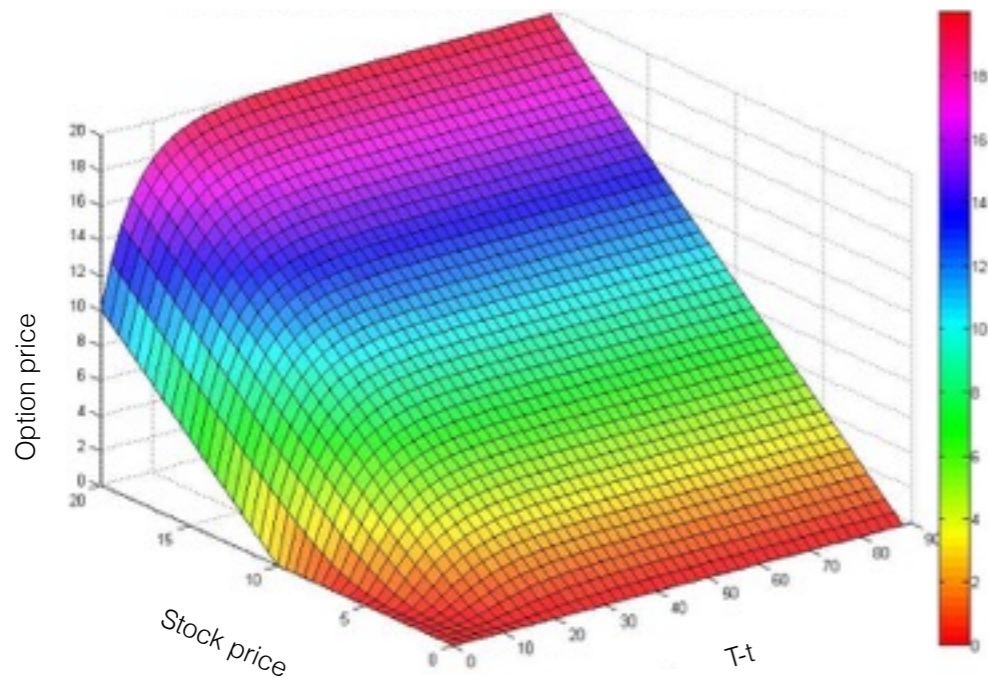
Where I come from: three application domains

1) DSL for Variational data assimilation



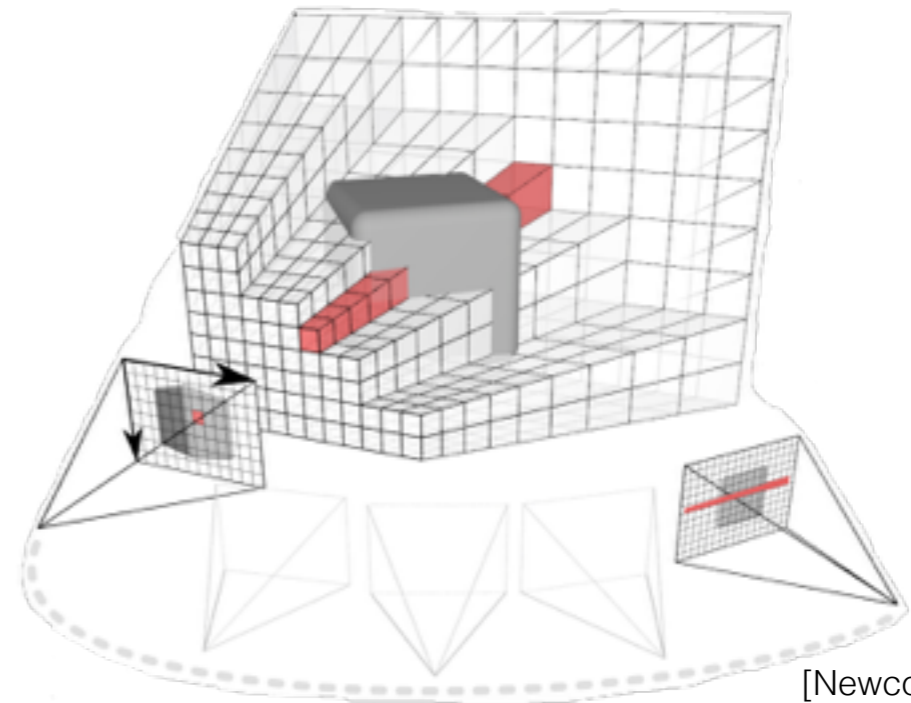
Ph.D. Computer/computational science, LOCEAN and CEDRIC labs

2) HPC computational finance



Permanent researcher, R&D at Murex S.A.S.
Murex Analytics (MACS) group

3) Computer vision

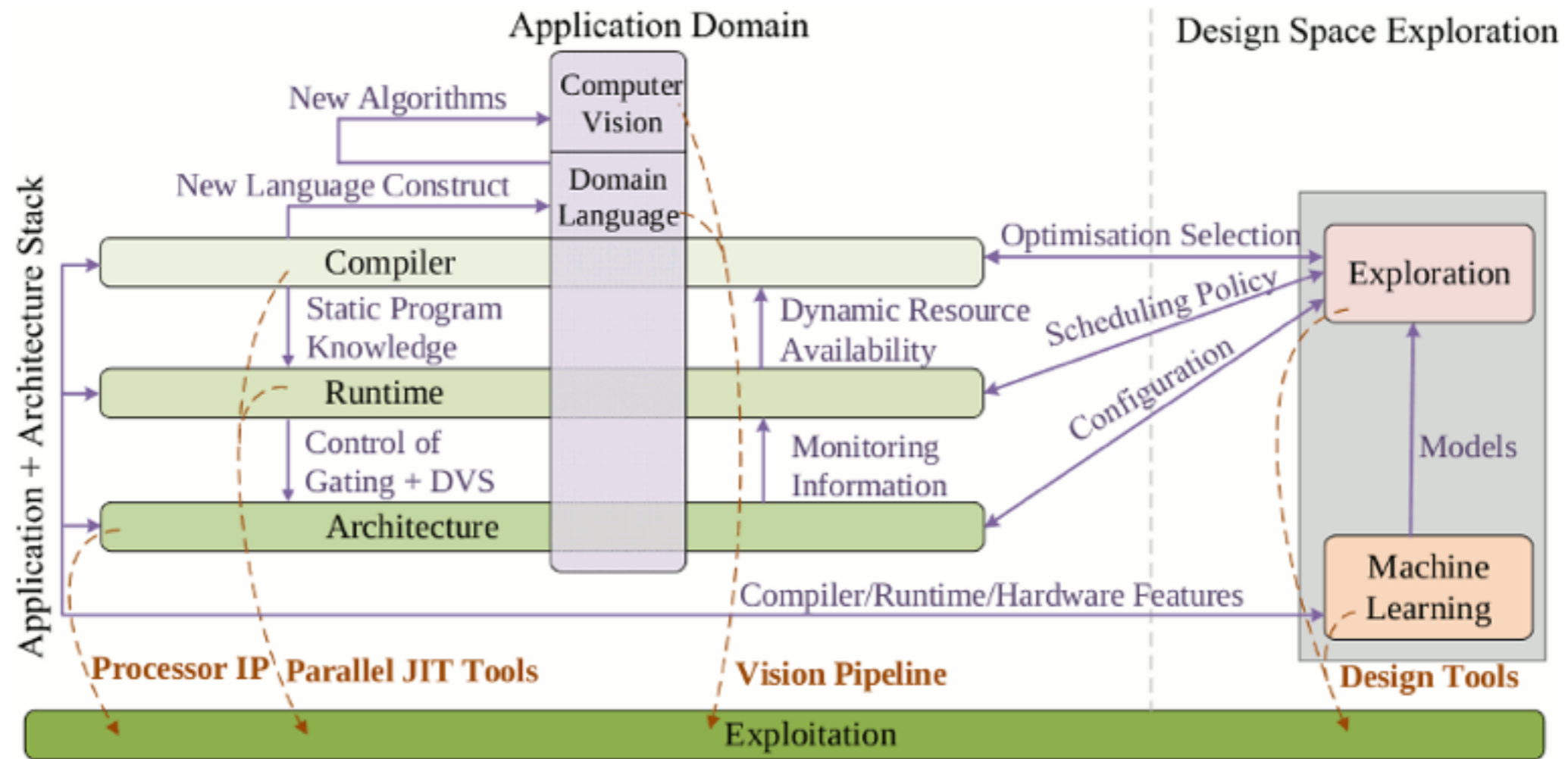


[Newcombe et al.,
ICCV 2011]

Research Associate, Imperial College London
Software Performance Optimisation group

PAMELA project

Panoramic Approach to the Many-core Landscape -
from application to end-device: a holistic approach

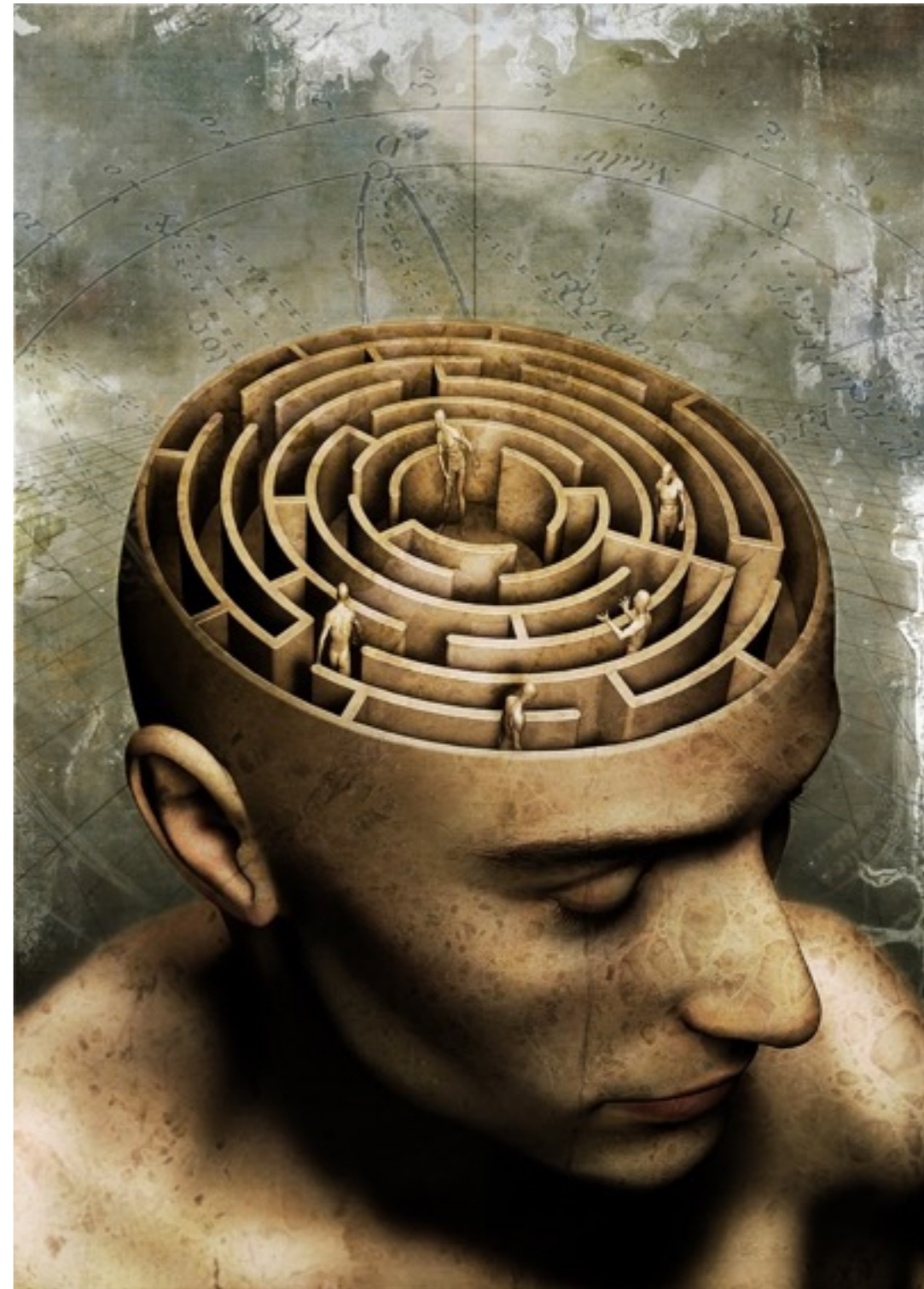


Main paper behind this talk

- Haifa, September 2016:
International Conference on
Parallel Architectures and Compilation Techniques (**PACT**)
- Paper title:
Integrating algorithmic parameters into benchmarking and
design space exploration in dense 3D scene understanding
- Authors:
B. Bodin, L. Nardi, Zia Zeeshan, H. Wagstaff, G. S. Shenoy, M.
Emani, J. Mawer, C. Kotselidis, A. Nisbet, M. Lujan, B. Franke, Paul H.
J. Kelly, M. O'Boyle

Outline

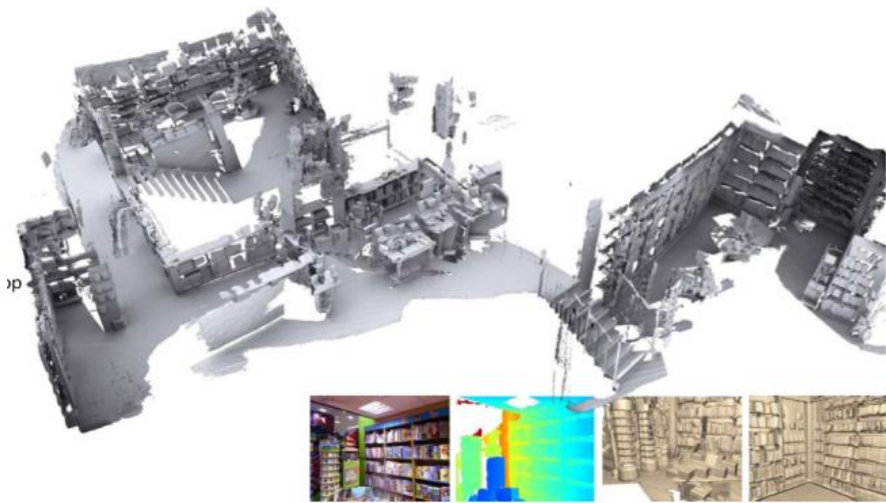
- **The SLAM application, a brief introduction**
- Benchmarking methodology
- Space exploration of algorithmic and implementation design choices



The three R's of vision: Spectrum of Computer Vision Research

Reconstruction

Scalable Kinect Fusion
(2013)

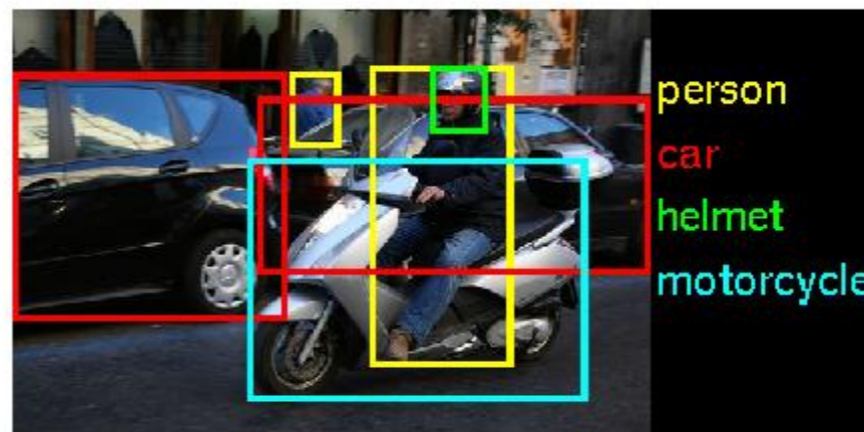
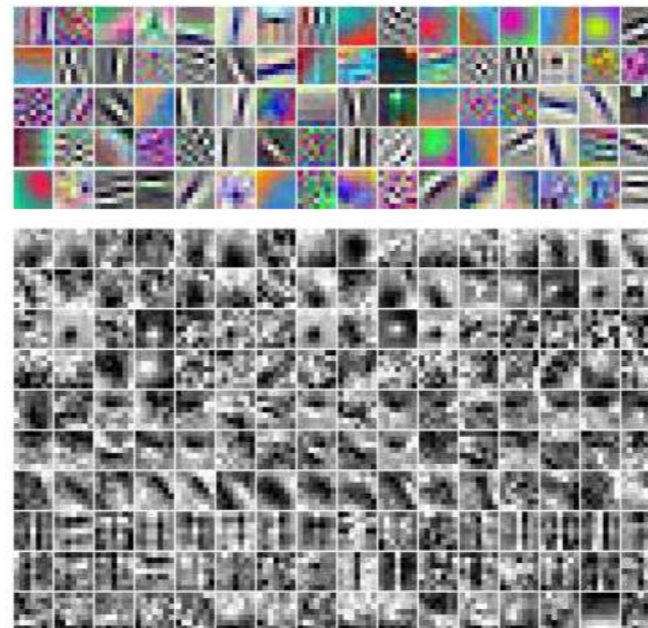


Building Rome on a
cloudless day (2010)



Recognition

Deep learning for scalable
Object class detection (2014)



Reorganisation or Grouping

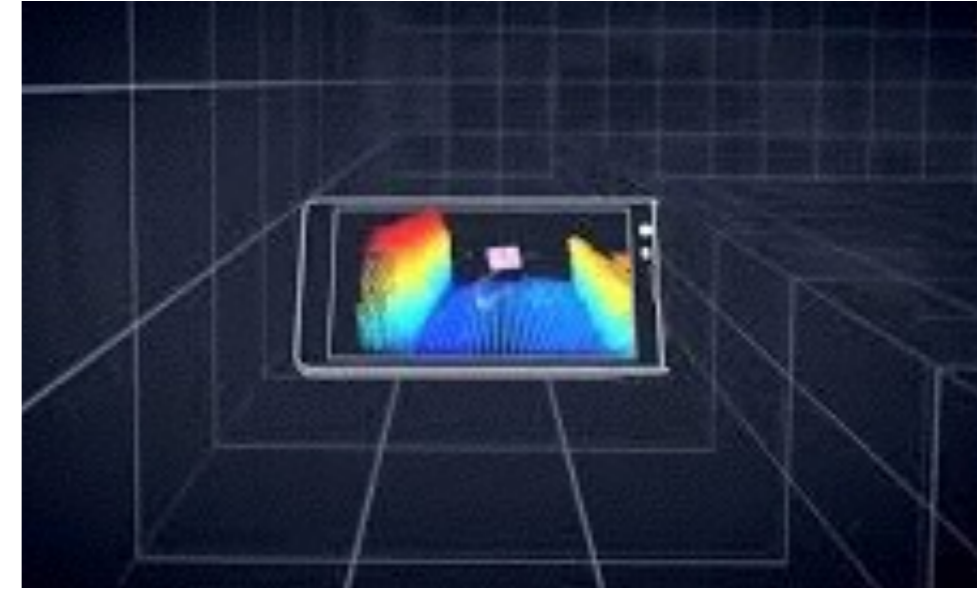
Contour detection
and segmentation (2011)



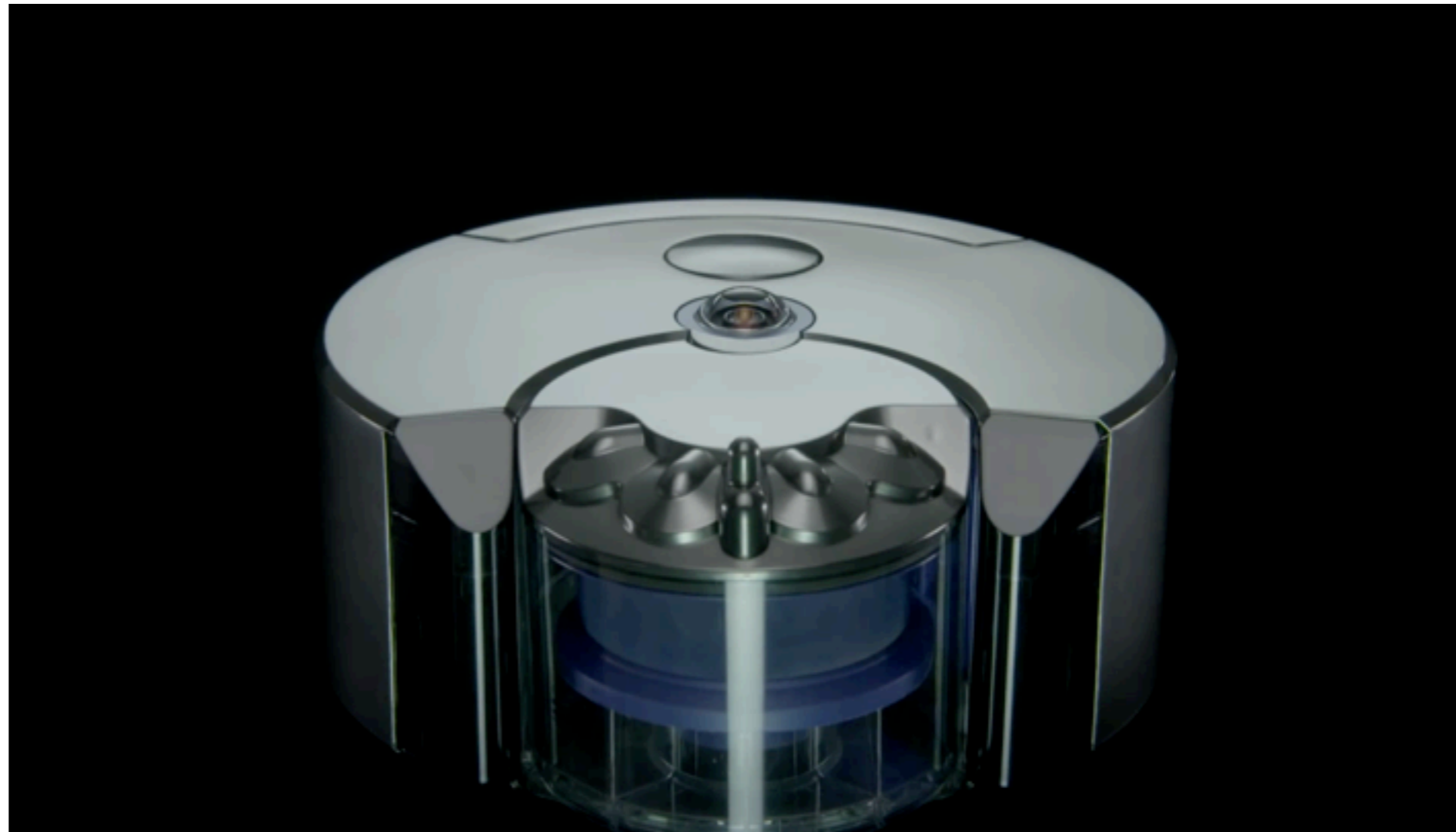
Simultaneous localisation and mapping (SLAM)

Build a coherent world representation and localise the camera in real-time

Sparse SLAM



Video:
[Dyson 360 Eye](#)



SIGGRAPH Talks 2011

KinectFusion:

Real-Time Dynamic 3D Surface
Reconstruction and Interaction

Shahram Izadi ¹, Richard Newcombe ², David Kim ^{1,3}, Otmar Hilliges ¹,
David Molyneaux ^{1,4}, Pushmeet Kohli ¹, Jamie Shotton ¹,
Steve Hodges ¹, Dustin Freeman ⁵, Andrew Davison ², Andrew Fitzgibbon ¹

¹ Microsoft Research Cambridge ² Imperial College London
³ Newcastle University ⁴ Lancaster University
⁵ University of Toronto

**Dense
SLAM**



Video: [KinectFusion](#)
[Newcombe et al. ISMAR 2011]

Simultaneous localisation and mapping (SLAM)

Build a coherent world representation and localise the camera in real-time

Dense SLAM

In this talk I will focus on two dense algorithms:

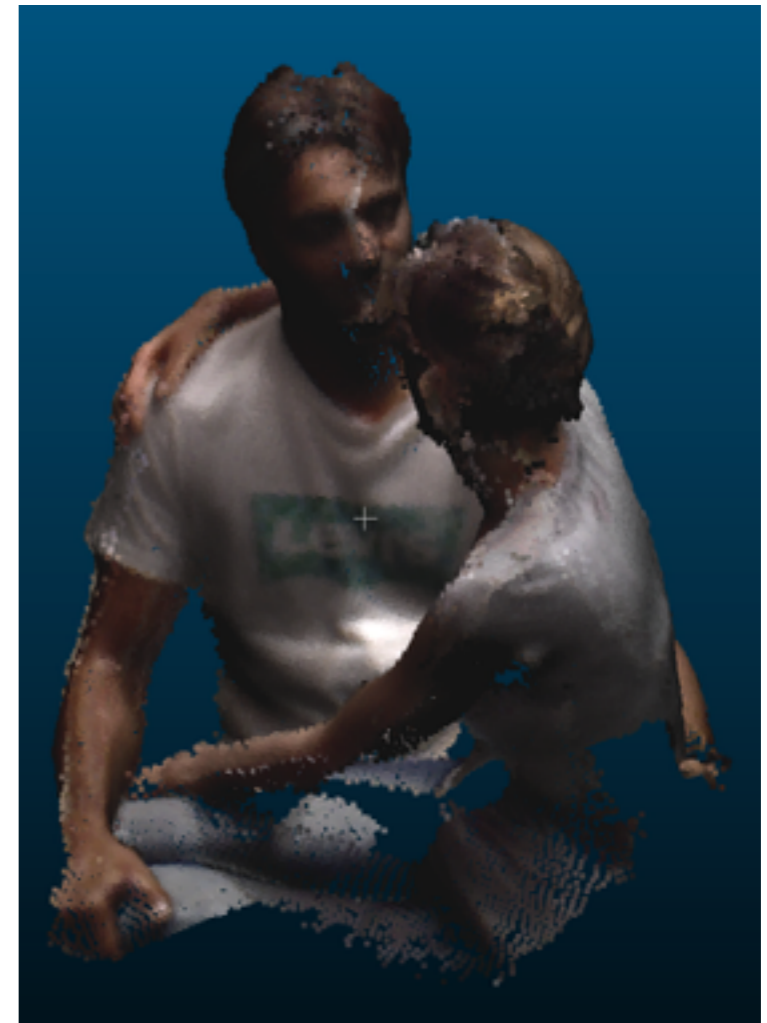
- KinectFusion [[Newcombe et al. ISMAR 2011](#)]
- ElasticFusion [[Whelan et al. RSS 2015](#)]

Applications, e.g.:

- Robotics
- Autonomous driving
- 3D printing
- Augmented reality
- Telepresence



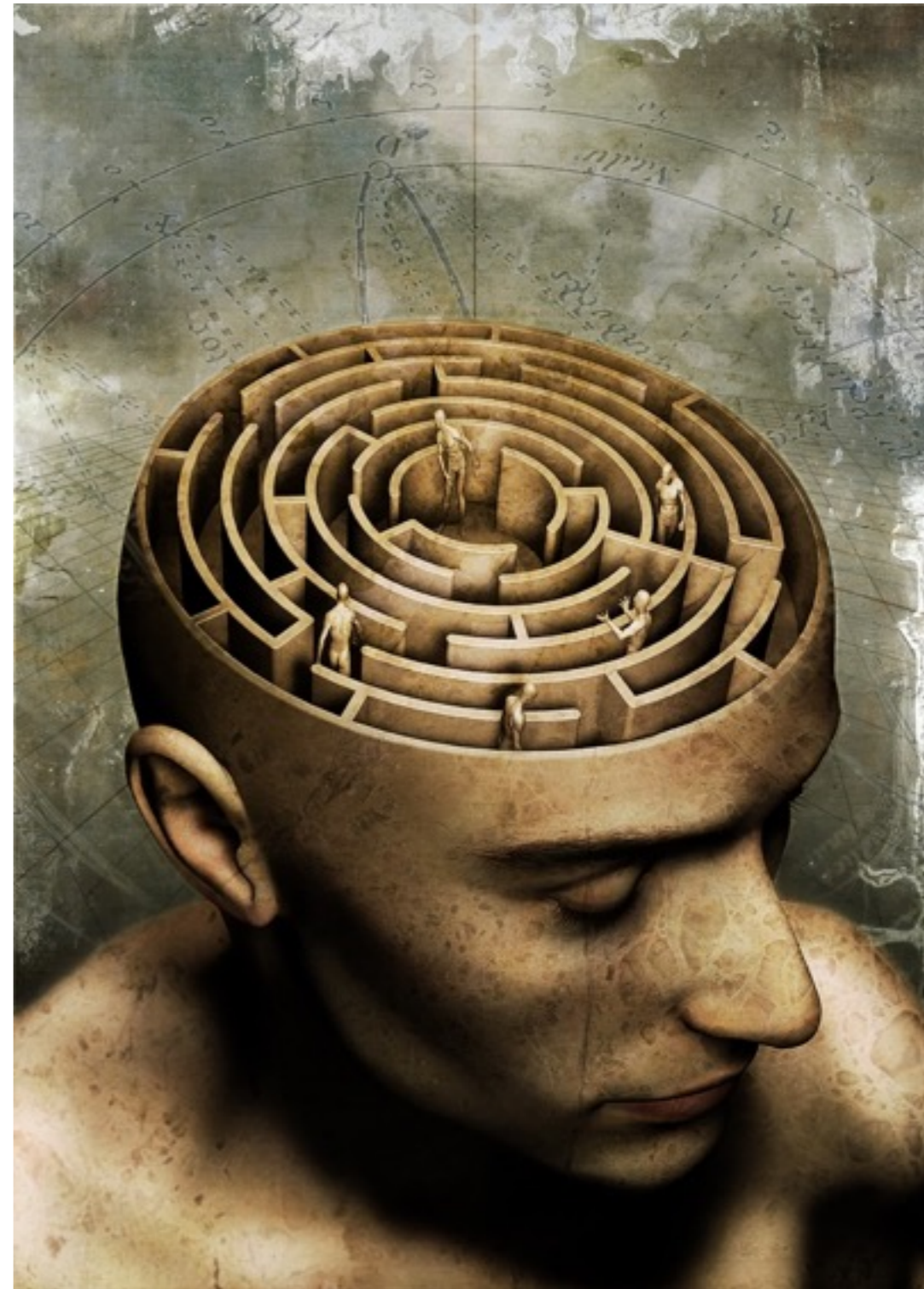
**Jesse Clayton (NVIDIA)
3D reconstruction**



**Daniele and Daniela
3D reconstruction**

Outline

- The SLAM application, a brief introduction
- **Benchmarking methodology**
- Space exploration of algorithmic and implementation design choices



What is “Performance”?

1. In several domains performance is **execution time**
2. In some domains performance is **accuracy**
3. What about **energy**?
4. But also memory consumption, temperature, robustness, etc.

A modern system evaluation considers multiple metrics:

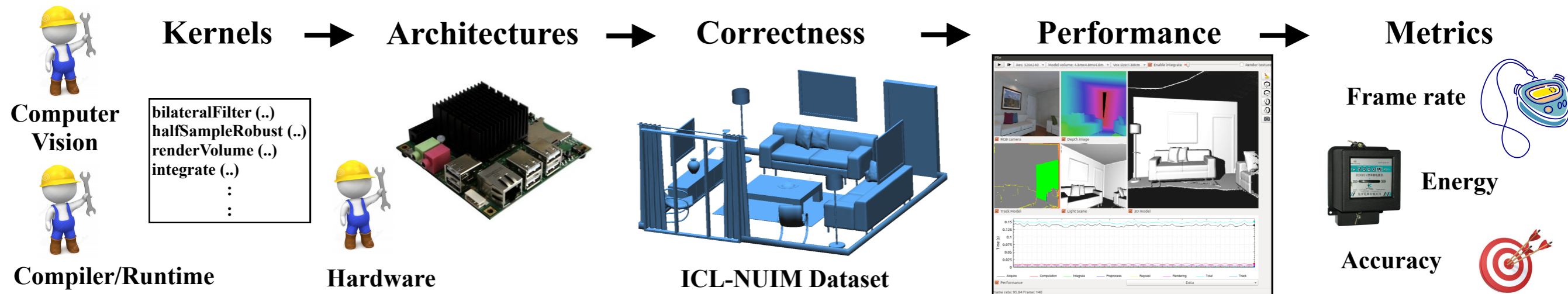
$$Performance = \begin{bmatrix} Runtime \\ Energy \\ Accuracy \end{bmatrix}$$

This defines a multi-objective optimisation problem: trade-off



Holistic approach to SLAM performance:

SLAMBench



A publicly-available benchmarking framework for quantitative, comparable and validatable experimental research to investigate trade-offs in performance, accuracy and energy consumption of a SLAM system

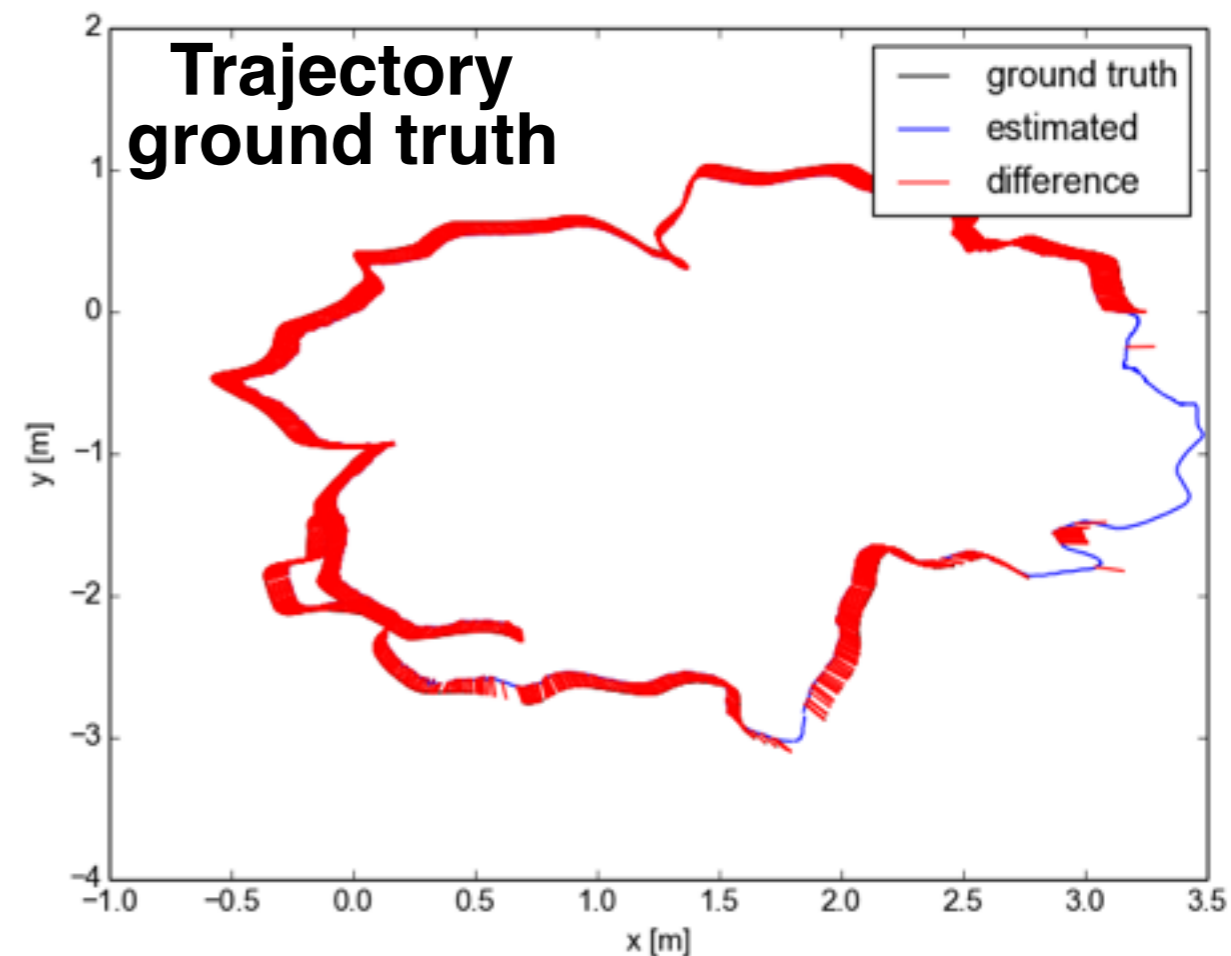
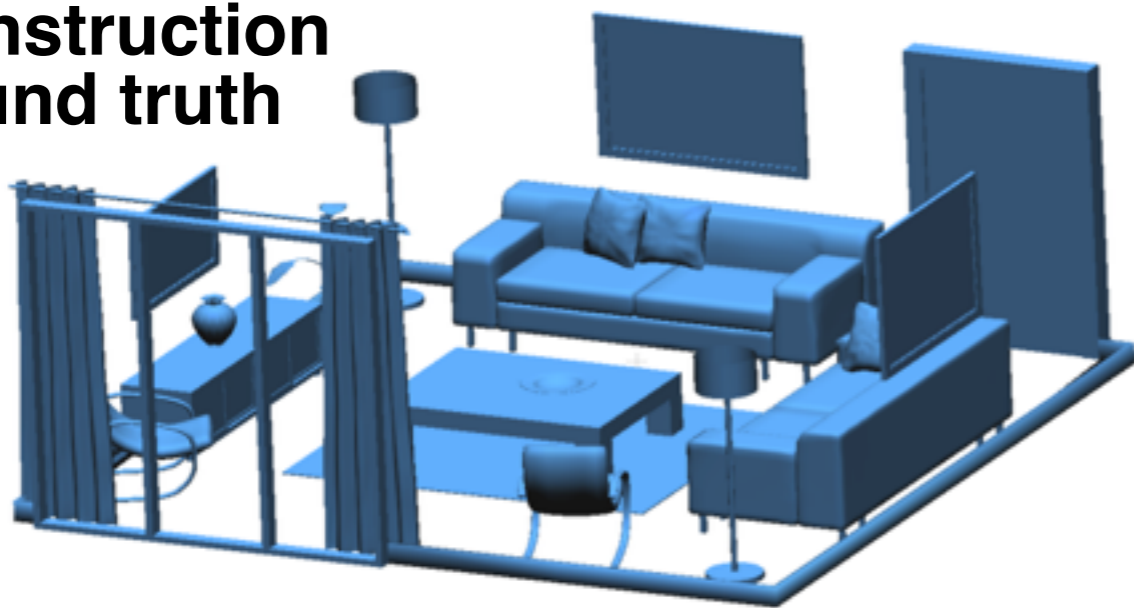
Introducing SLAMBench, a performance and accuracy benchmarking methodology for SLAM (ICRA 2015)



ICL-NUIM dataset

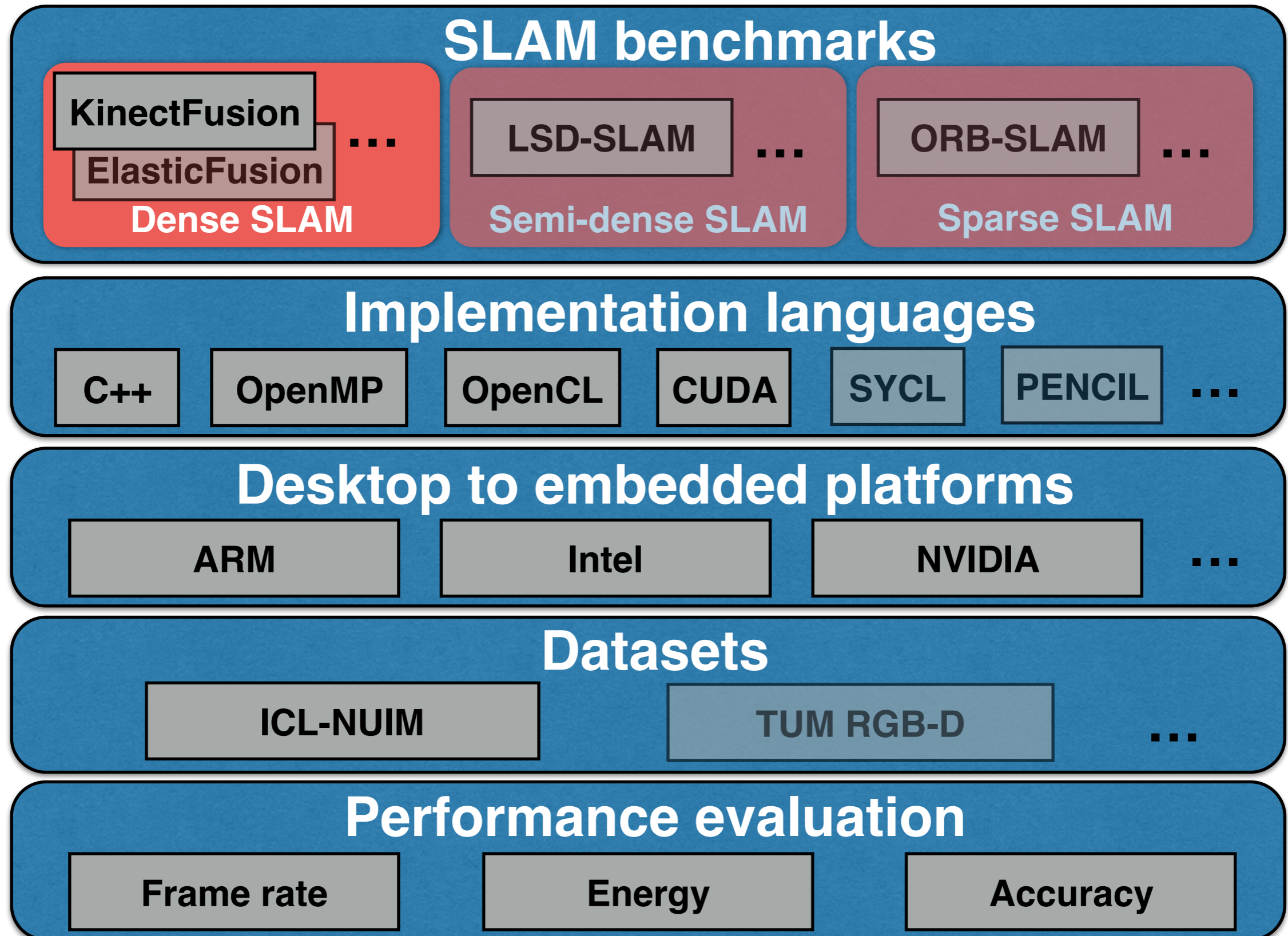


Reconstruction ground truth

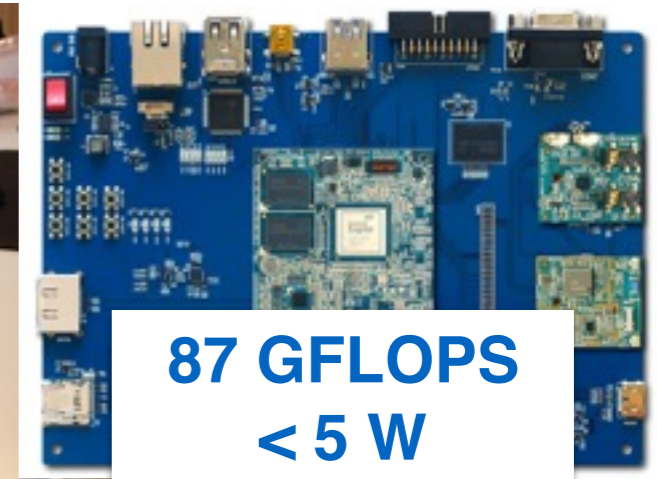
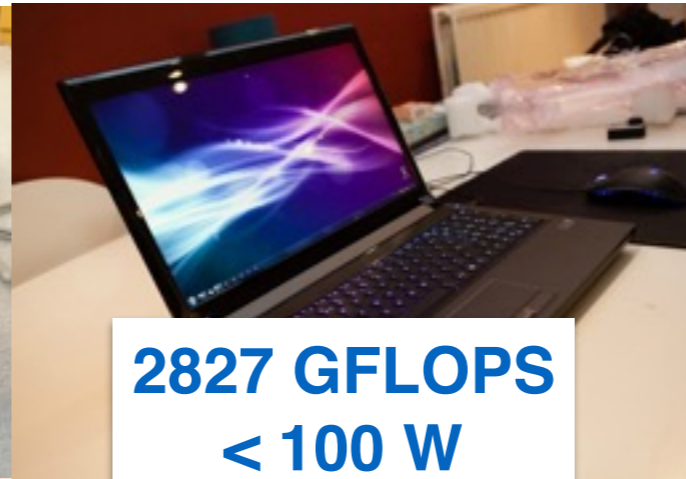


- ICL-NUIM synthetic dataset [Handa et al. 2014]
- 880 RGB-D frames at 30 FPS
- Absolute trajectory error (ATE) based on ground truth

SLAMBench framework



Machines	TITAN	GTX870M	TK1	ODROID	Arndale
CPU	Intel	Intel	ARM	ARM	ARM
CPU name	i7 Haswell	i7 Haswell	NVIDIA 4-Plus-1	Exynos 5422	Exynos 5250
CPU GFLOPS	448	307	74	80	27
CPU cores	4	4	4 + 1	4 + 4	2
GPU	NVIDIA	NVIDIA	NVIDIA	ARM	ARM
GPU name	TITAN	GTX 870M	Tegra K1	Mali-T628	Mali-T604
GPU GFLOPS	4500	2520	330	60 + 30	60



Platforms



“Performance” on SLAMBench

- Runtime/energy/accuracy measurements
- Accuracy provided via absolute trajectory error (ATE)



Machine	CPU	CPU name	CPU GFLOPS	CPU cores	GPU	GPU name	GPU GFLOPS	TDP Watts
Hardkernel ODROID-XU3	ARM A15 + A7	Exynos 5422	80	4 + 4	ARM	Mali-T628	60 + 30	10

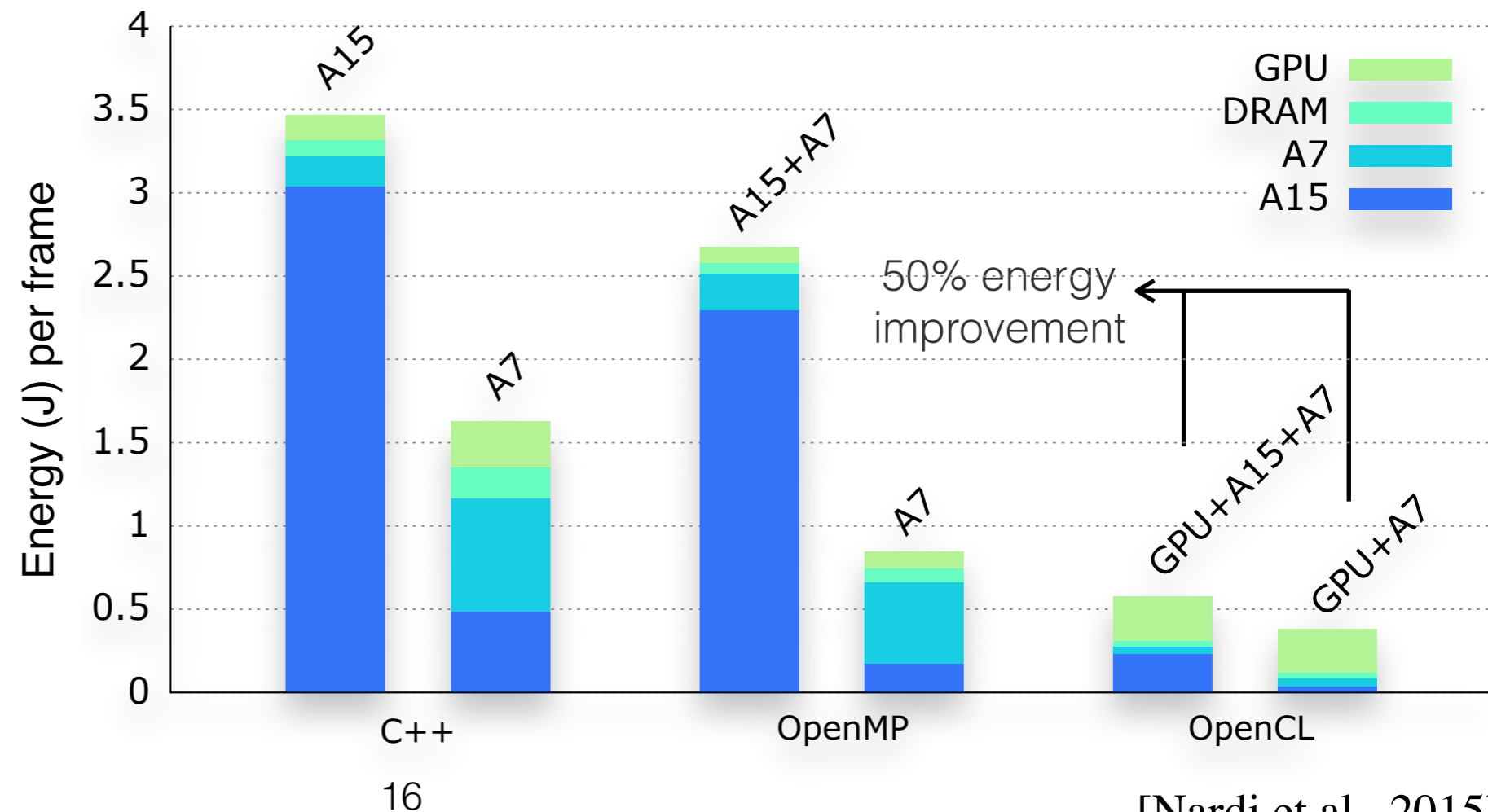
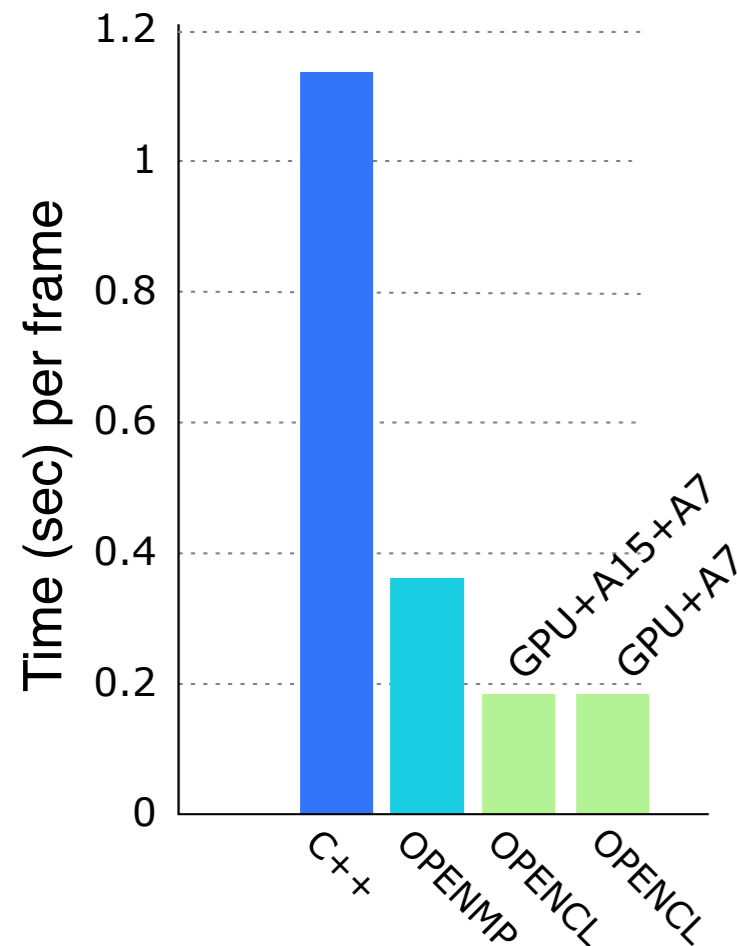
“Performance” on SLAMBench

- Runtime/energy/accuracy measurements
- Accuracy provided via absolute trajectory error (ATE)

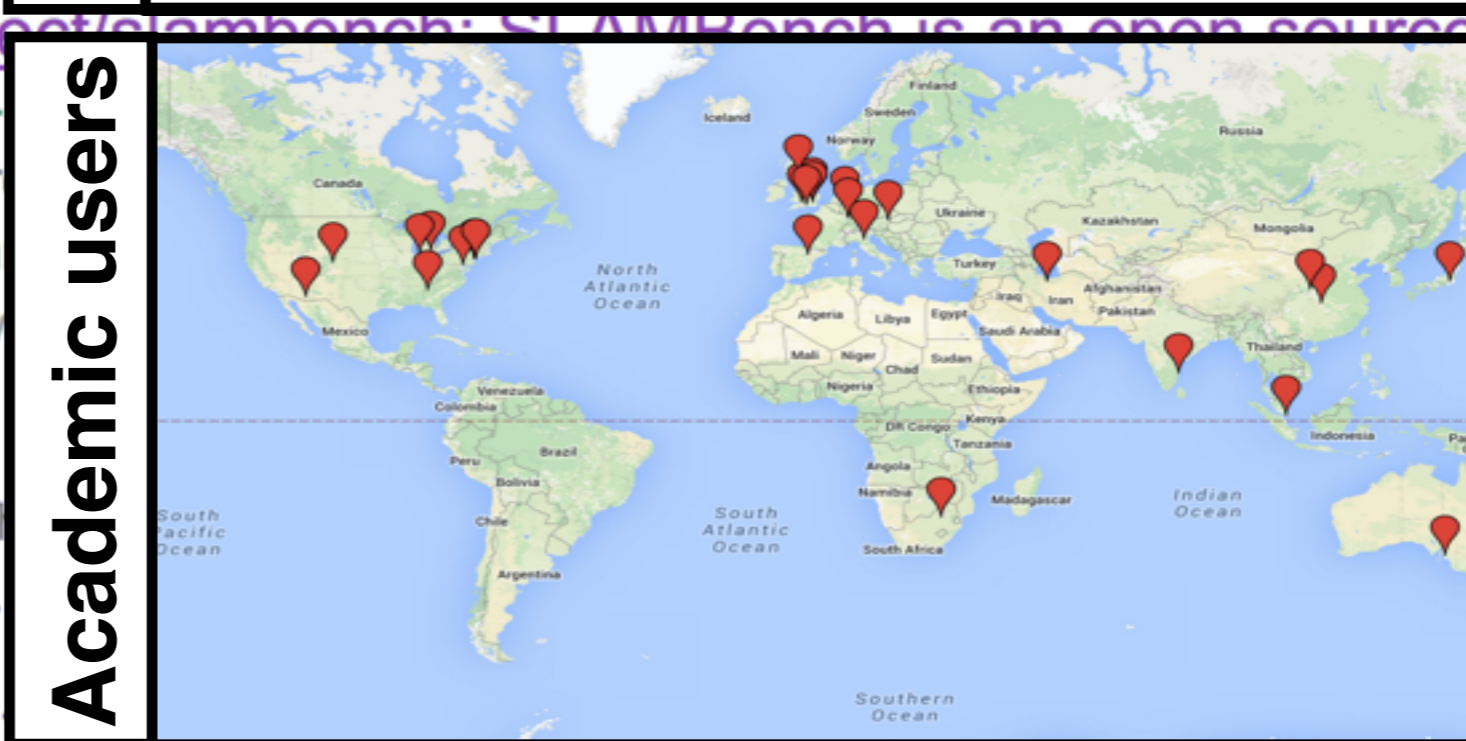


ATE in cm	
C++	2.06
OpenMP	2.06
OpenCL	2.01

Machine	CPU	CPU name	CPU GFLOPS	CPU cores	GPU	GPU name	GPU GFLOPS	TDP Watts
Hardkernel ODROID-XU3	ARM A15 + A7	Exynos 5422	80	4 + 4	ARM	Mali-T628	60 + 30	10



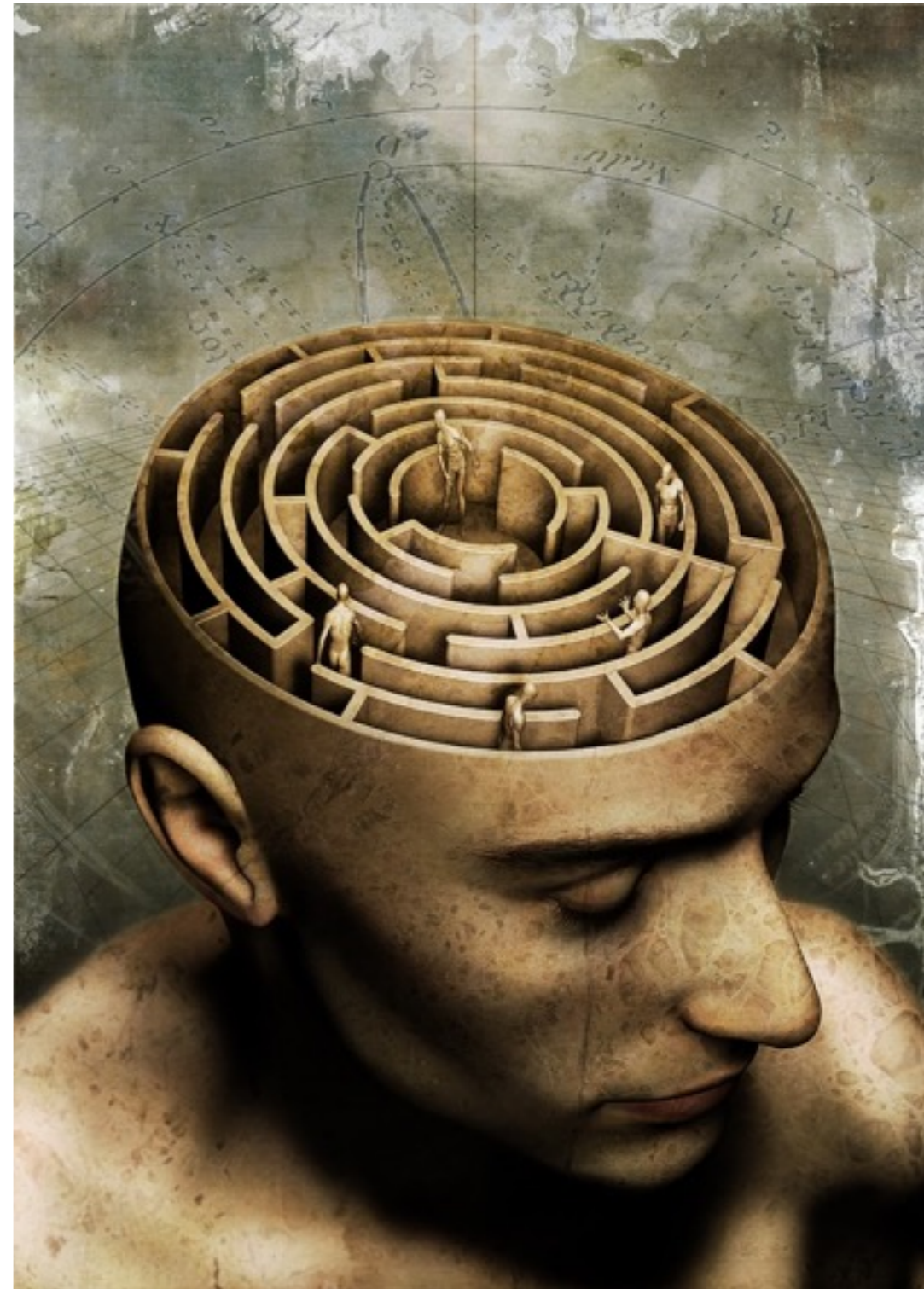
Google search results for "slambench github". The search bar shows "slambench github" and the results include "About 1,240 results (0.64 sec)". The top result is "GitHub - pamela-project/slambench: SLAMBench is an open source tool designed to...". Below the title is a "README.md" link and the text "slambench - SLAMBench source tool designed to...". A link for "More results from github" is also visible.



Publicly released 13/11/2014
 (1400+ downloads)

Outline

- The SLAM application, a brief introduction
- Benchmarking methodology
- **Space exploration of algorithmic and implementation design choices**



What is the optimisation space?

Configuration parameters:

Space 1

1. Algorithmic:
 - Application-specific parameters
 - Minimisation methods
 - Early exit condition values

What is the optimisation space?

Configuration parameters:

Space 1	<ol style="list-style-type: none">1. Algorithmic:<ul style="list-style-type: none">• Application-specific parameters• Minimisation methods• Early exit condition values
Space 2	<ol style="list-style-type: none">2. Compilation:<ul style="list-style-type: none">• opencl-params: -cl-mad-enable, -cl-fast-relaxed-math, etc.• LLVM flags: O1, O2, O3, vectorize-slp-aggressive, etc.• Local work group size: 16/32/64/96/112/128/256• Vectorisation: width (1/2/4/8), direction (x/y)• Thread coarsening: factor (1/2/4/8/16/32), stride (1/2/4/8/16/32), dimension (x/y)

What is the optimisation space?

Configuration parameters:

Space 1	<p>1. Algorithmic:</p> <ul style="list-style-type: none"> • Application-specific parameters • Minimisation methods • Early exit condition values
Space 2	<p>2. Compilation:</p> <ul style="list-style-type: none"> • opencl-params: -cl-mad-enable, -cl-fast-relaxed-math, etc. • LLVM flags: O1, O2, O3, vectorize-slp-aggressive, etc. • Local work group size: 16/32/64/96/112/128/256 • Vectorisation: width (1/2/4/8), direction (x/y) • Thread coarsening: factor (1/2/4/8/16/32), stride (1/2/4/8/16/32), dimension (x/y)
Space 3	<p>3. Architecture:</p> <ul style="list-style-type: none"> • GPU frequency: 177/266/350/420/480/543/600/DVFS • # of active big cores: 0/1/2/3/4 • # of active LITTLE cores: 1/2/3/4

What is the optimisation space?

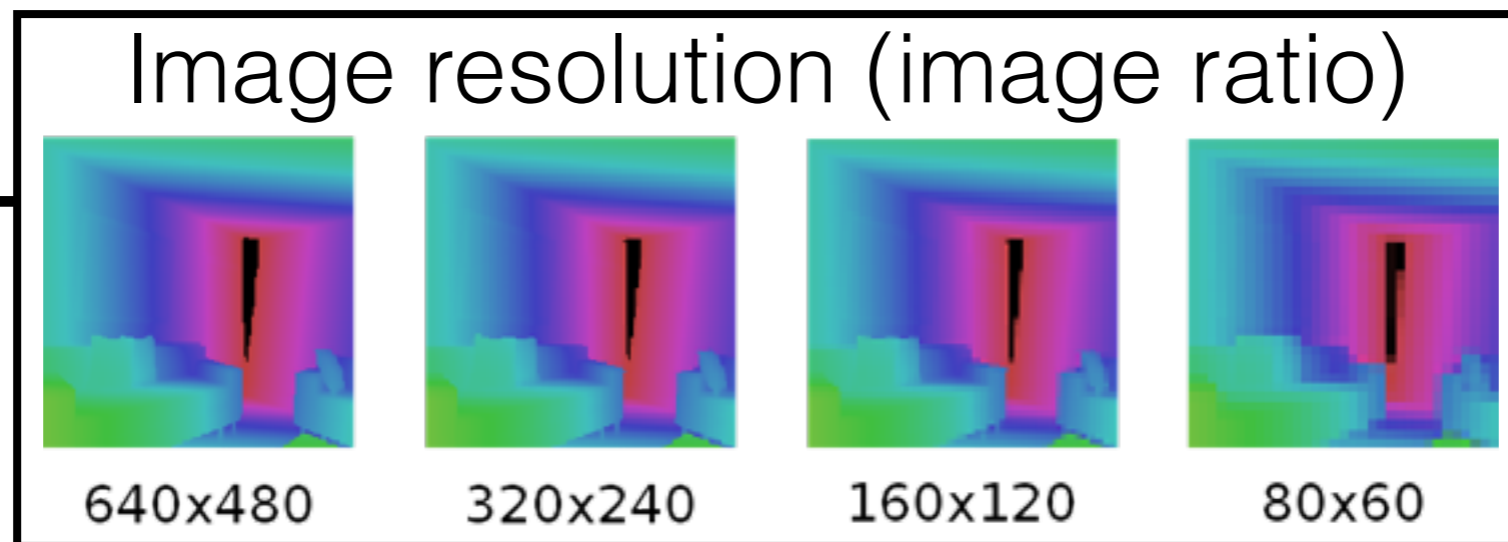
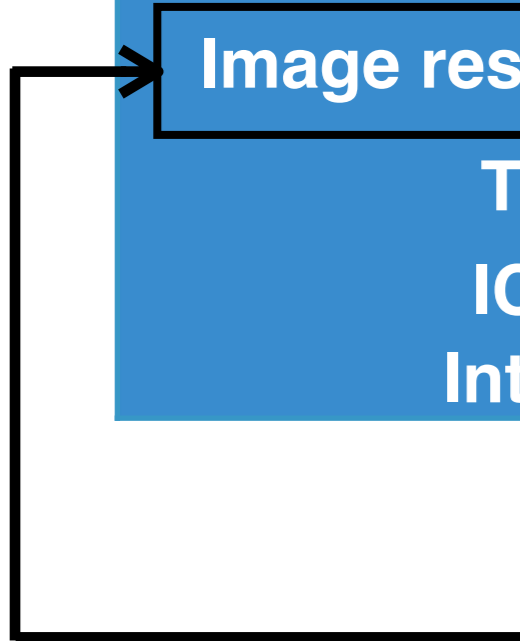
Configuration parameters:

Co-design space	Space 1	<ol style="list-style-type: none"> Algorithmic: <ul style="list-style-type: none"> Application-specific parameters Minimisation methods Early exit condition values
	Space 2	<ol style="list-style-type: none"> Compilation: <ul style="list-style-type: none"> opencl-params: -cl-mad-enable, -cl-fast-relaxed-math, etc. LLVM flags: O1, O2, O3, vectorize-slp-aggressive, etc. Local work group size: 16/32/64/96/112/128/256 Vectorisation: width (1/2/4/8), direction (x/y) Thread coarsening: factor (1/2/4/8/16/32), stride (1/2/4/8/16/32), dimension (x/y)
	Space 3	<ol style="list-style-type: none"> Architecture: <ul style="list-style-type: none"> GPU frequency: 177/266/350/420/480/543/600/DVFS # of active big cores: 0/1/2/3/4 # of active LITTLE cores: 1/2/3/4

Warning: huge spaces, impossible to run exhaustively

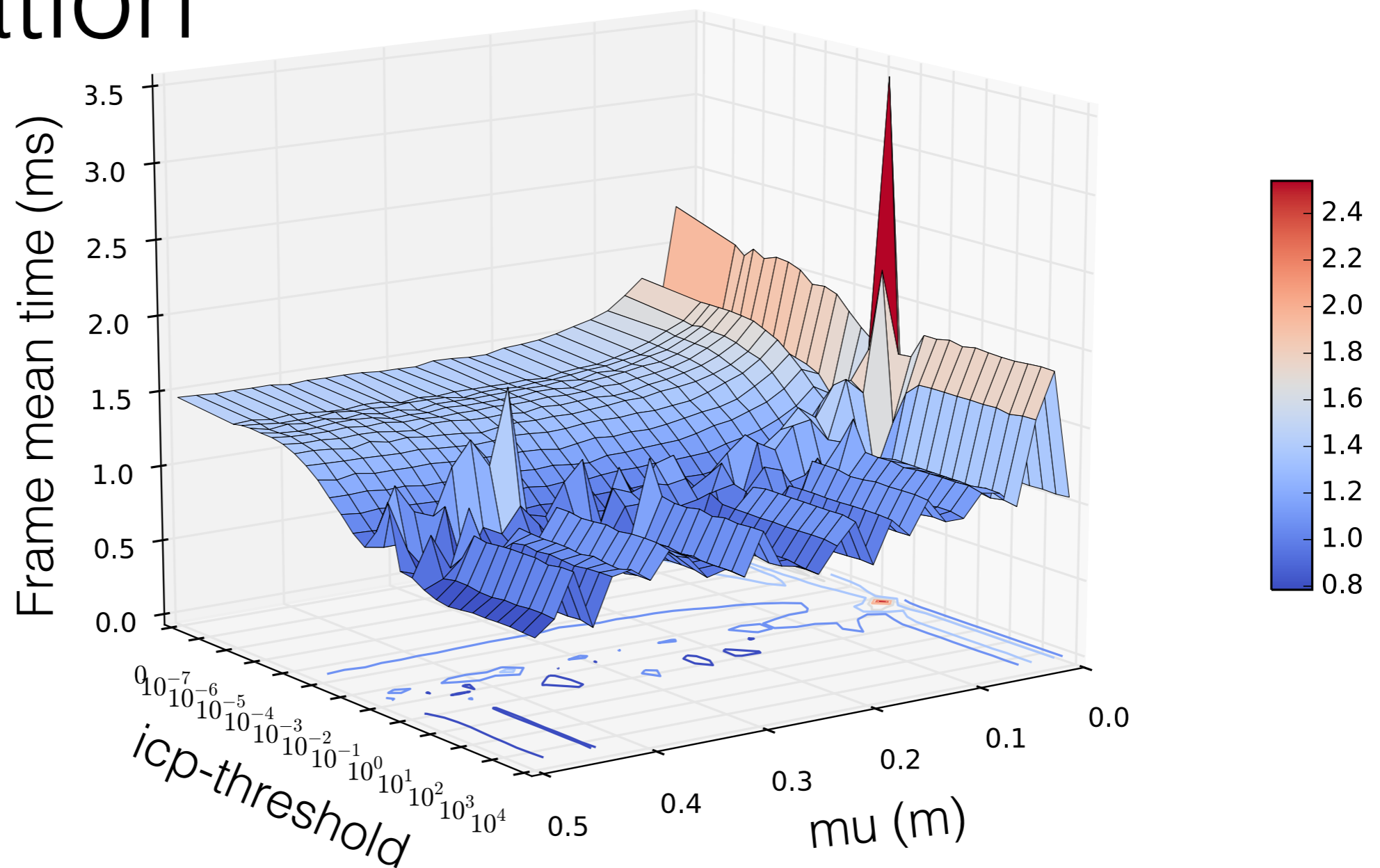
KinectFusion algorithmic features

Features	Ranges
Volume resolution	64x64x64, 128x128x128, 256x256x256, 512x512x512
μ distance	0 .. 0.5
Pyramid level iterations (3 levels)	0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11
Image resolution (image ratio)	1, 2, 4, 8
Tracking rate	1, 2, 3, 4, 5
ICP threshold	10^{-6} .. 10^2
Integration rate	1 .. 30



Different algorithmic features for ElasticFusion

Motivation

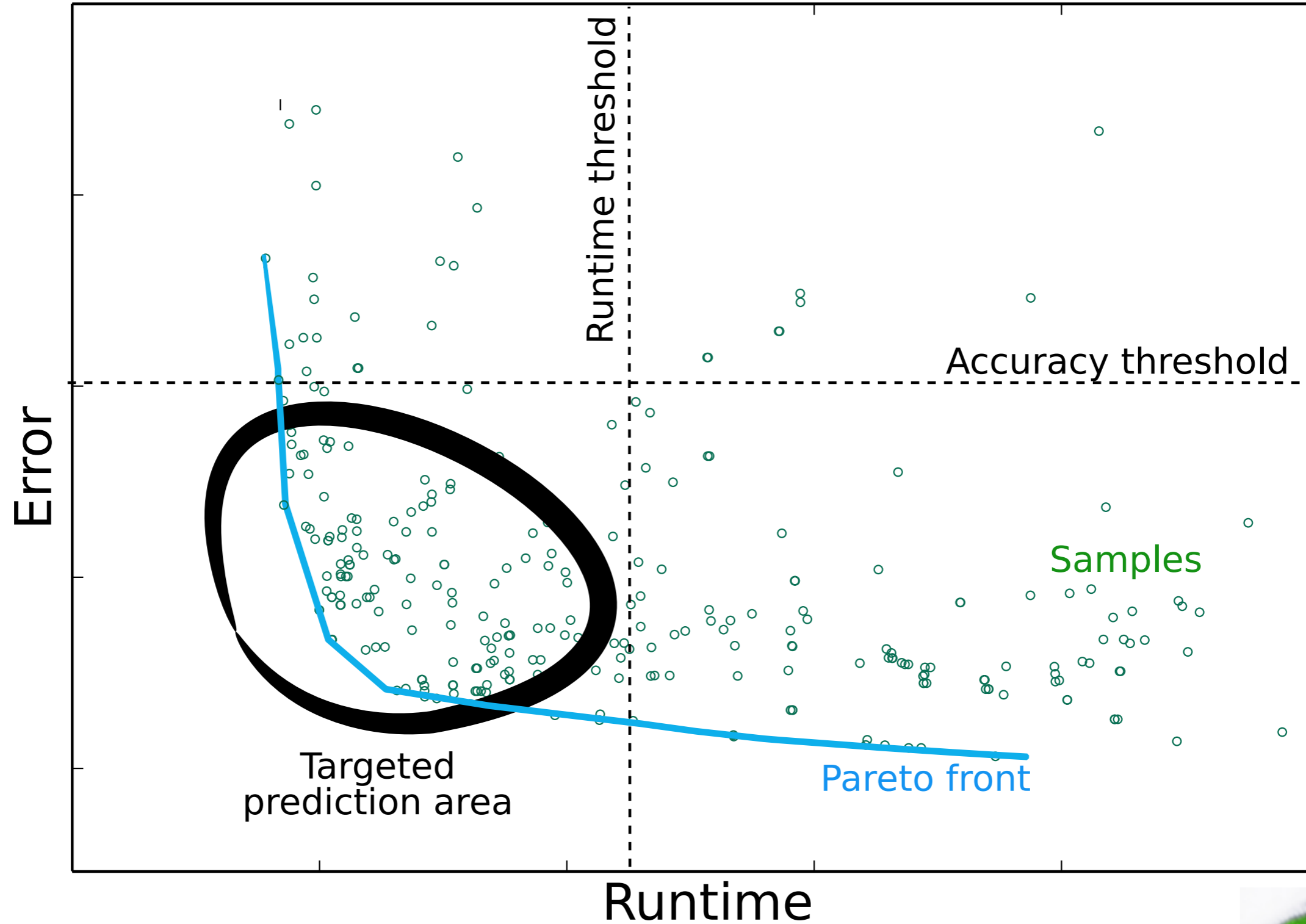


- KinectFusion runtime response surface: non-linear, multi-modal and non-smooth
- Optimal **algorithm configurability** enables better performance and better accuracy of the computation

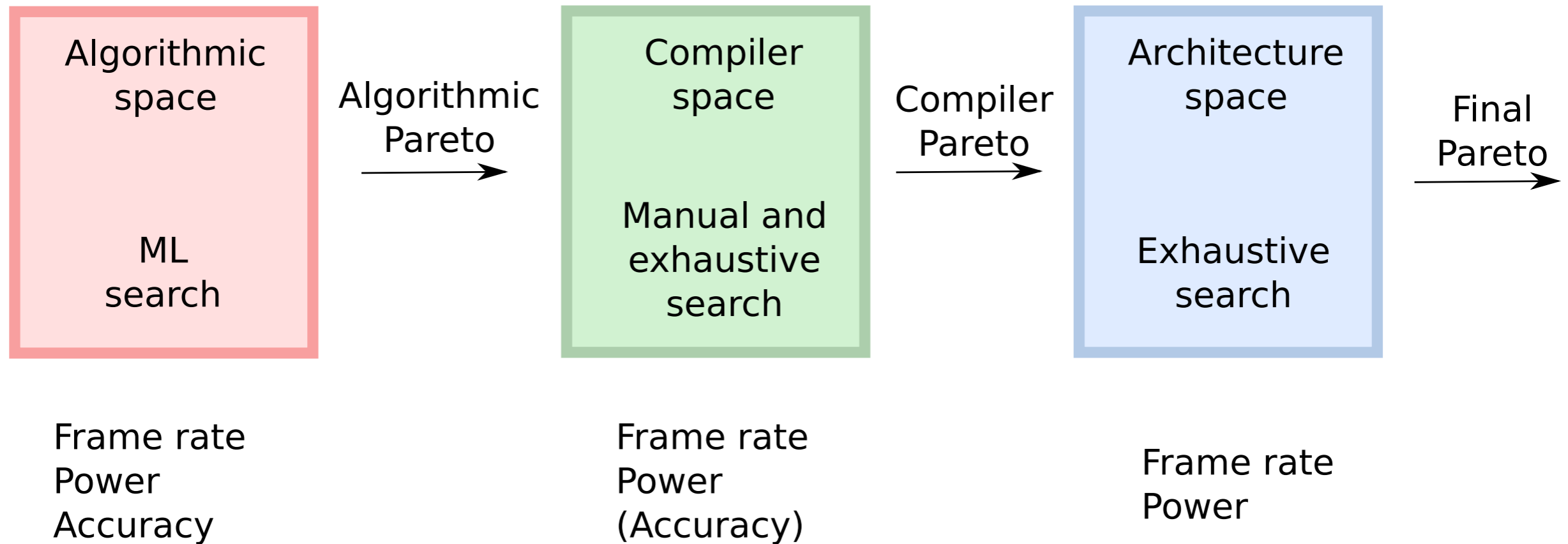
Integrating Algorithmic Parameters into Benchmarking and Design Space Exploration in 3D Scene Understanding (PACT 2016)



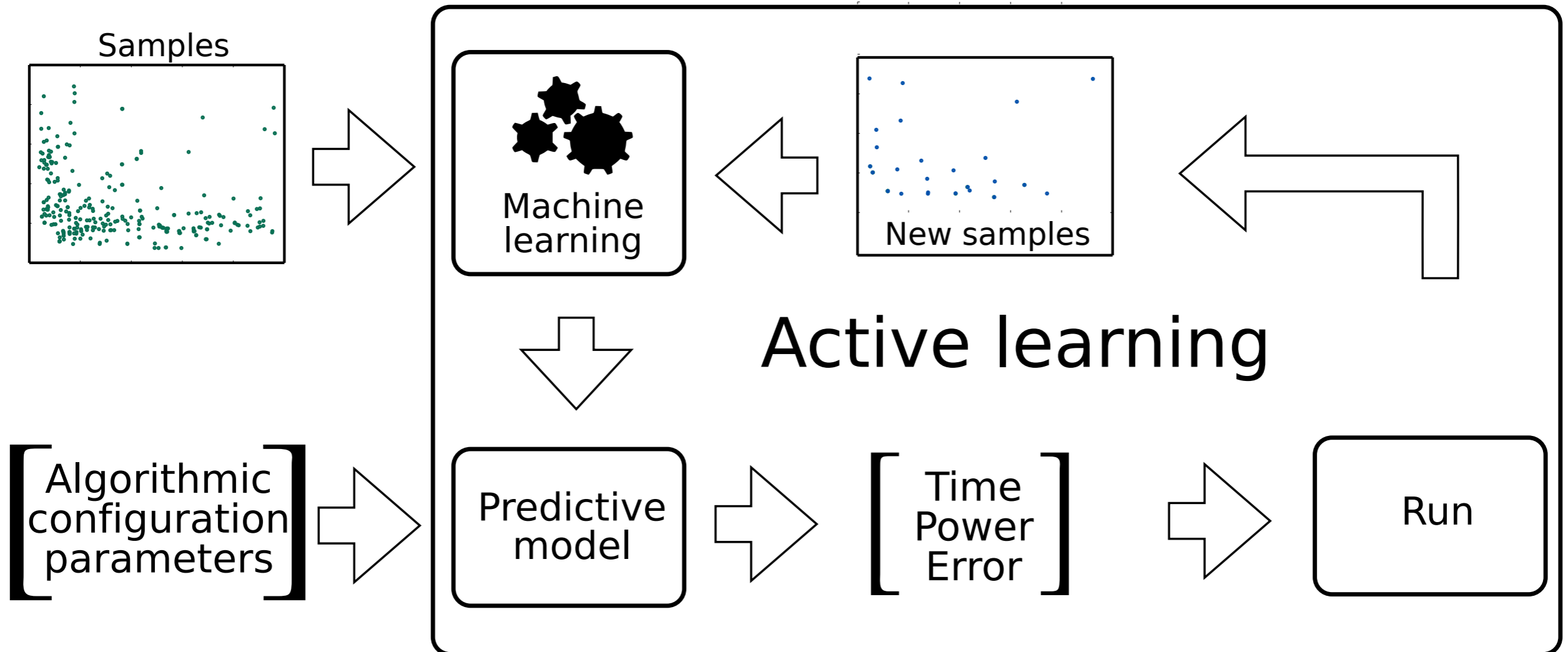
Exploration goal



Incremental exploration approach

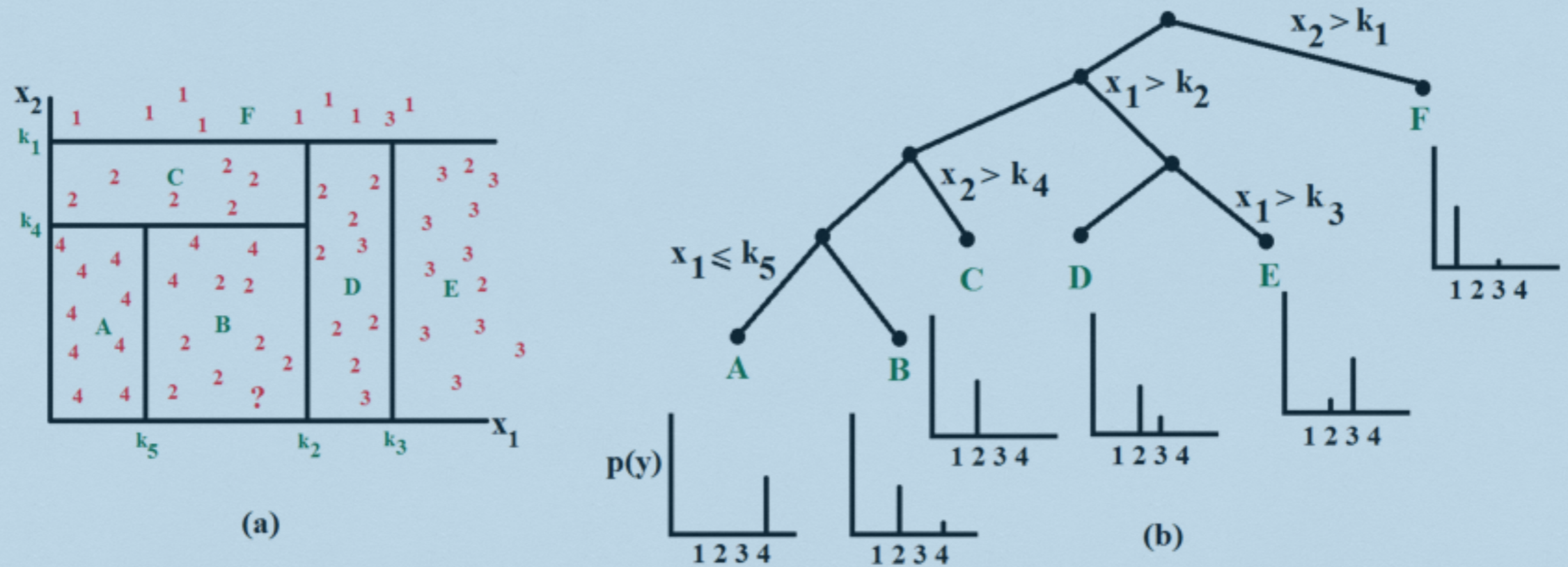


Algo design-space exploration (DSE)

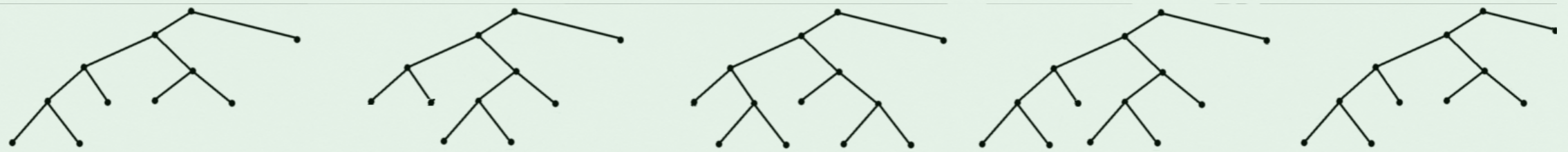


Machine learning methods used

Decision Tree

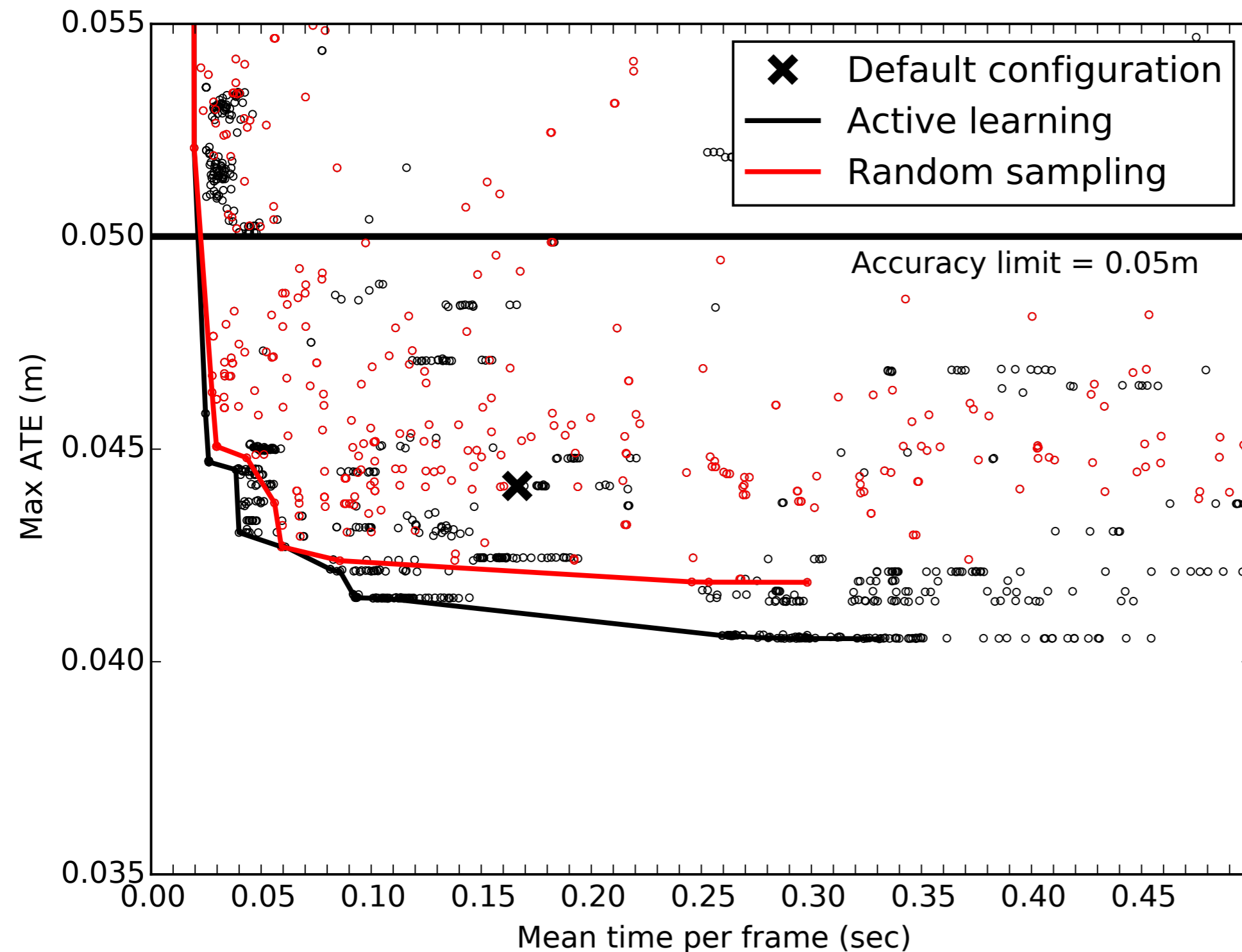


Random Forest



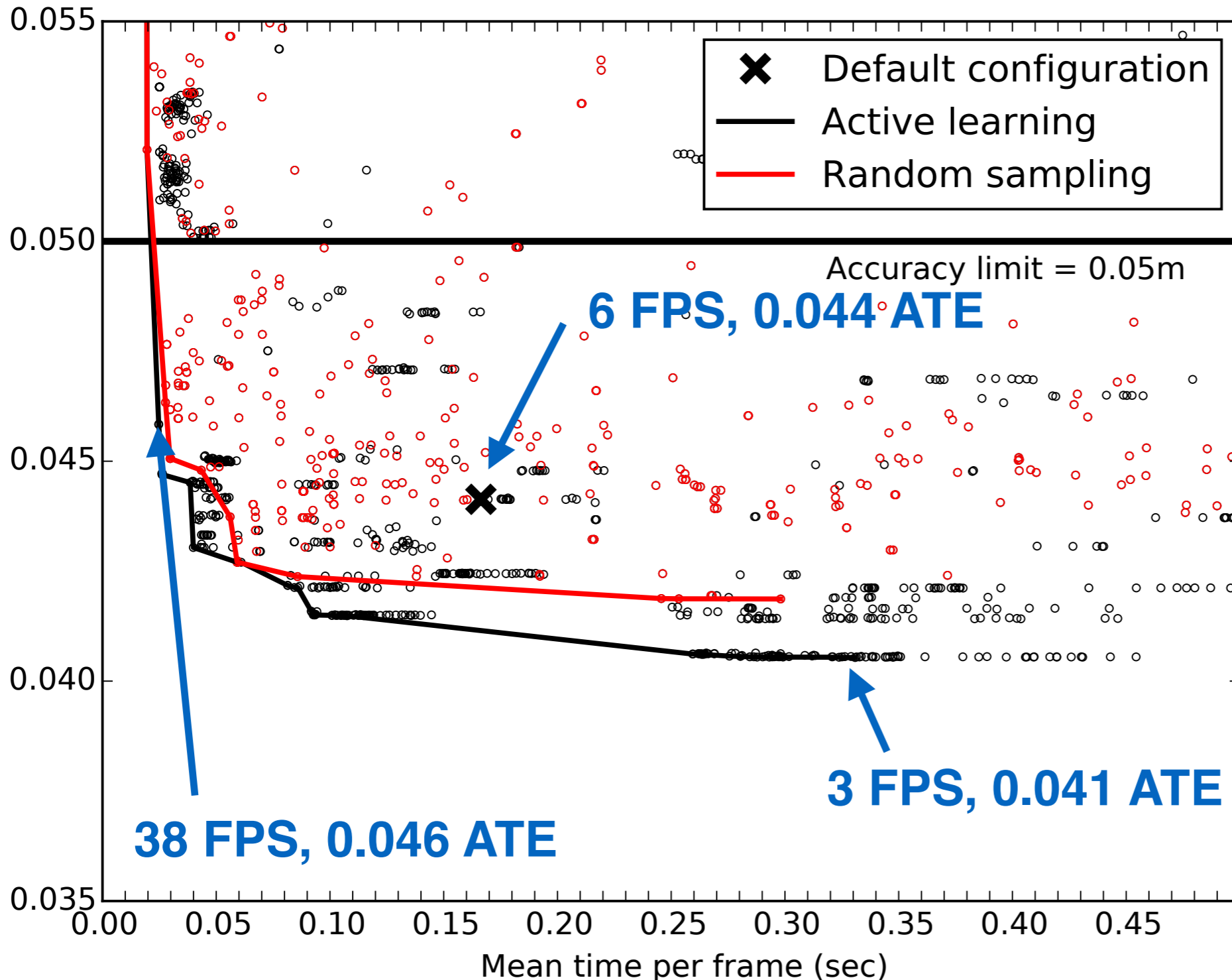
Results **KinectFusion** algo DSE error/runtime

Machine	Type	CPU	CPU name	CPU cores	GPU	GPU name
Hardkernel ODROID-XU3	Embedded	ARM A15 + A7	Exynos 5422	4 + 4	ARM	Mali-T628

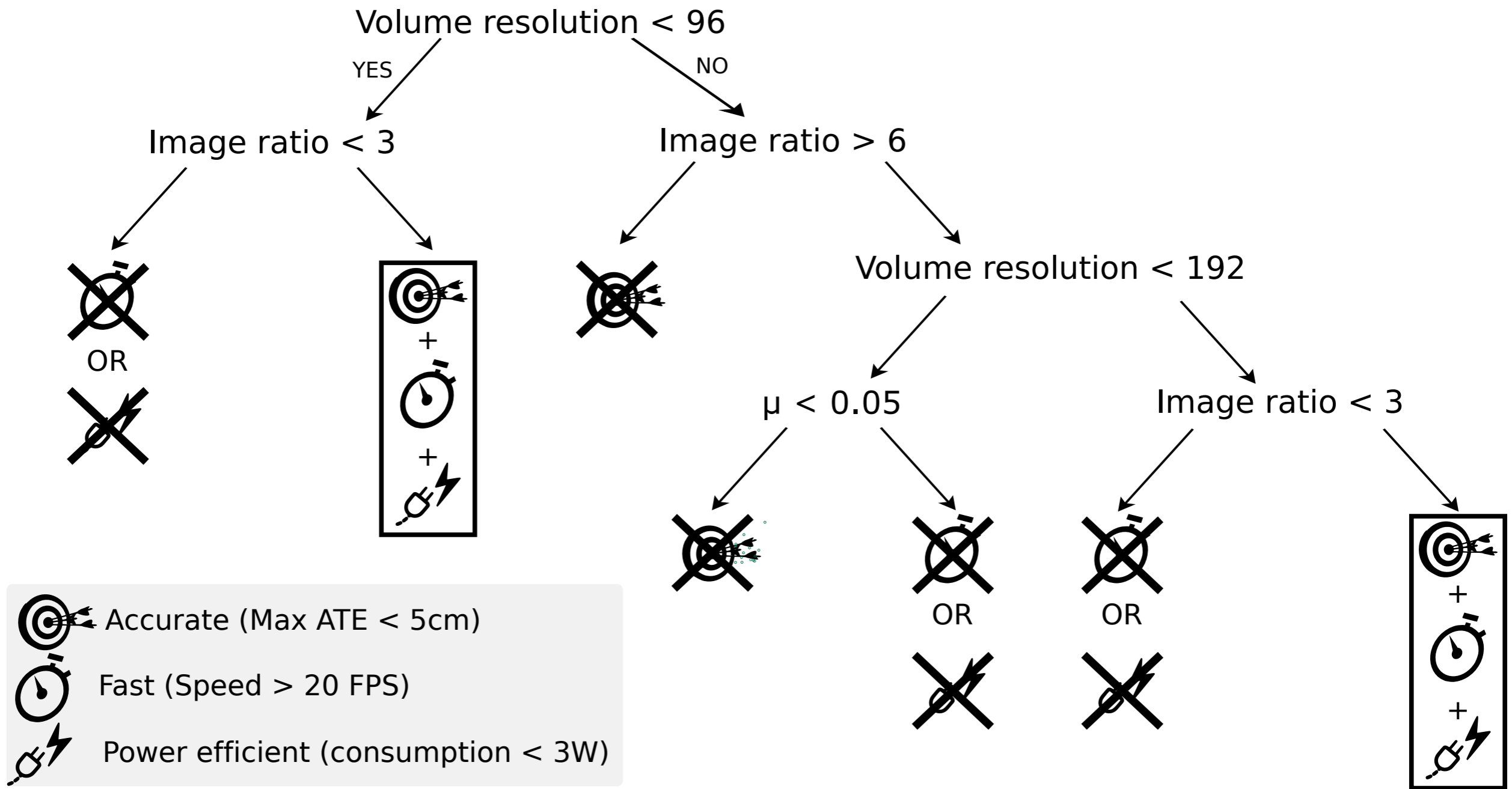


Results KinectFusion algo DSE error/runtime

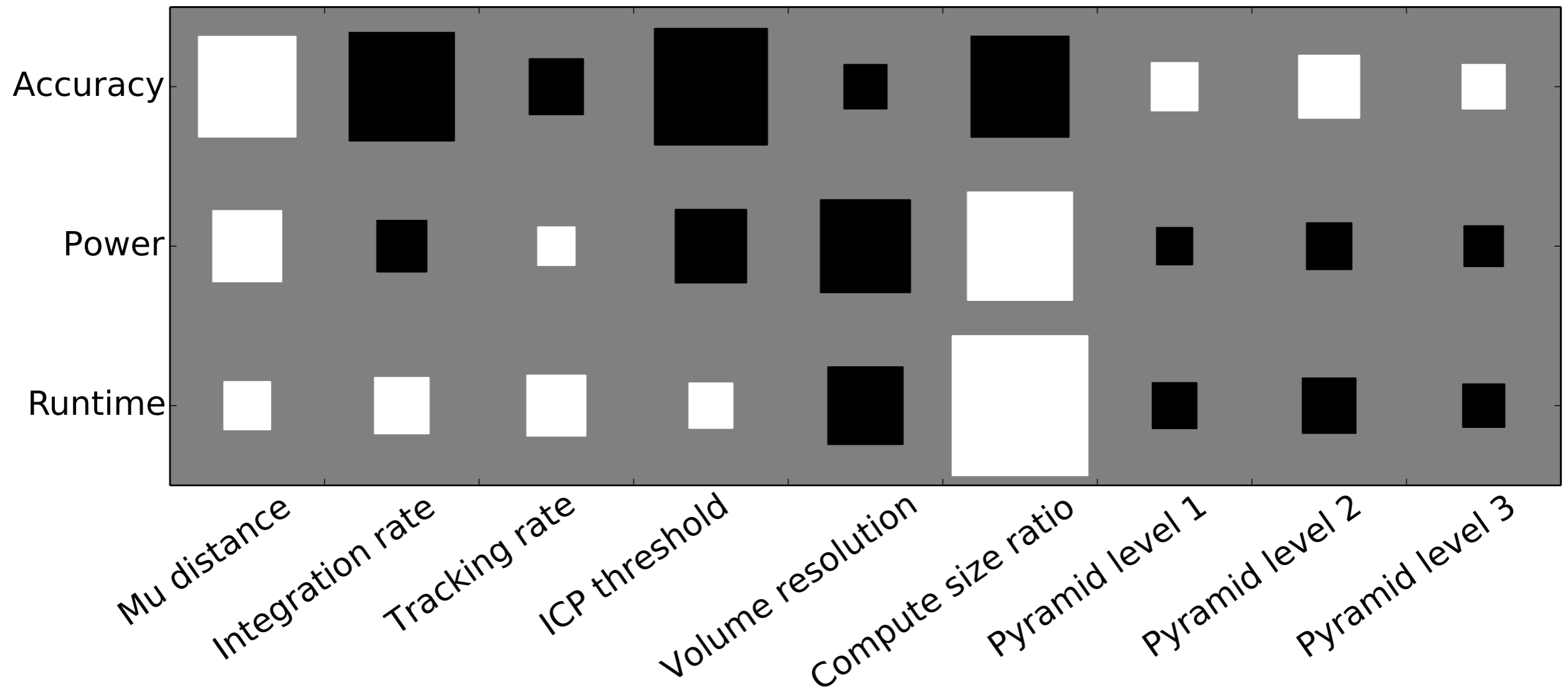
Machine	Type	CPU	CPU name	CPU cores	GPU	GPU name
Hardkernel ODROID-XU3	Embedded	ARM A15 + A7	Exynos 5422	4 + 4	ARM	Mali-T628



Predominant algorithmic features



Hinton correlation diagram on algorithmic features



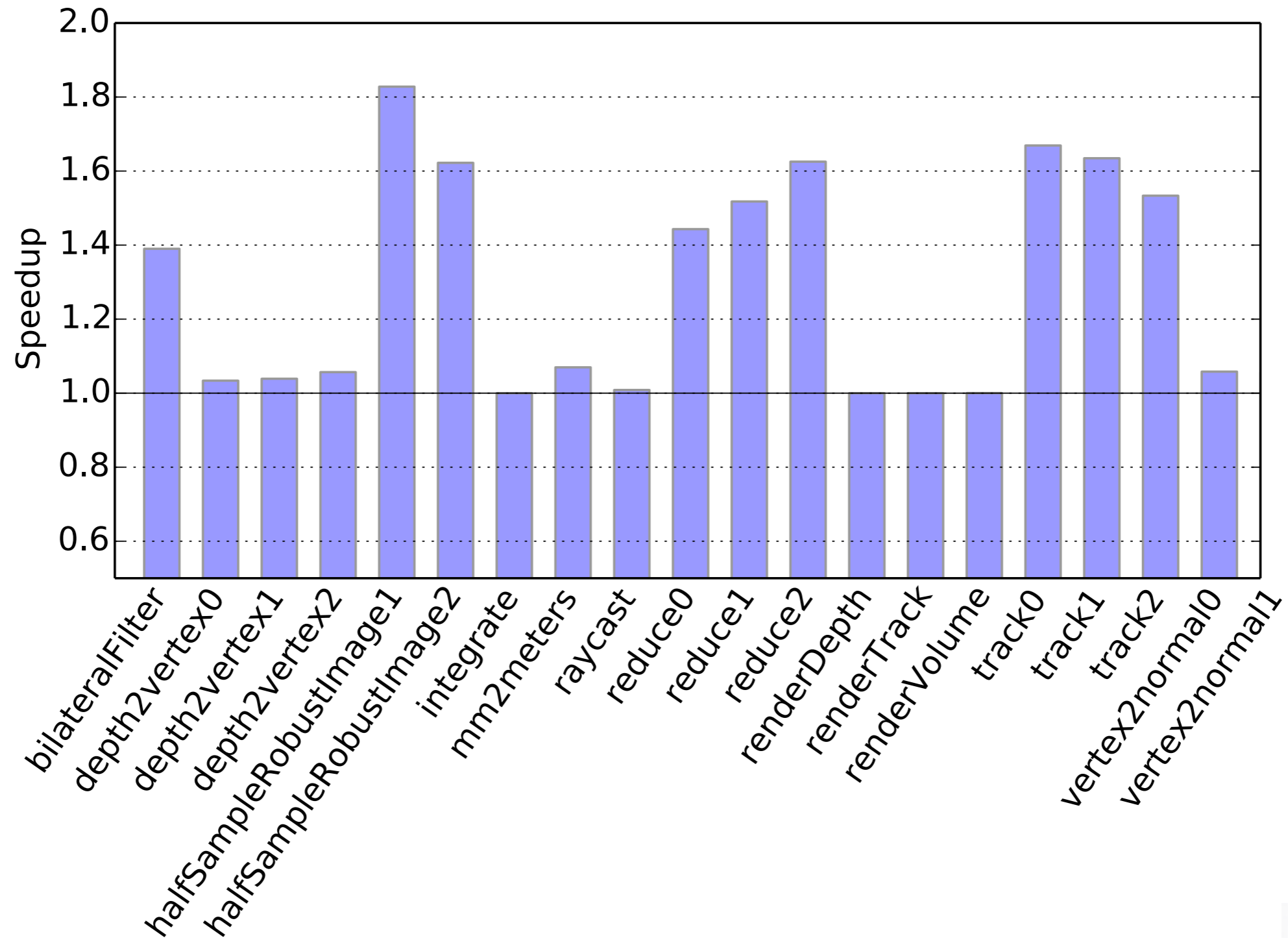
- Impact of algorithmic parameters (x-axis) on the performance metrics (y-axis) on the ODROID-XU3 platform.
- Bigger squares indicate a higher correlation.

White square: a parameter which when increased improves the metric

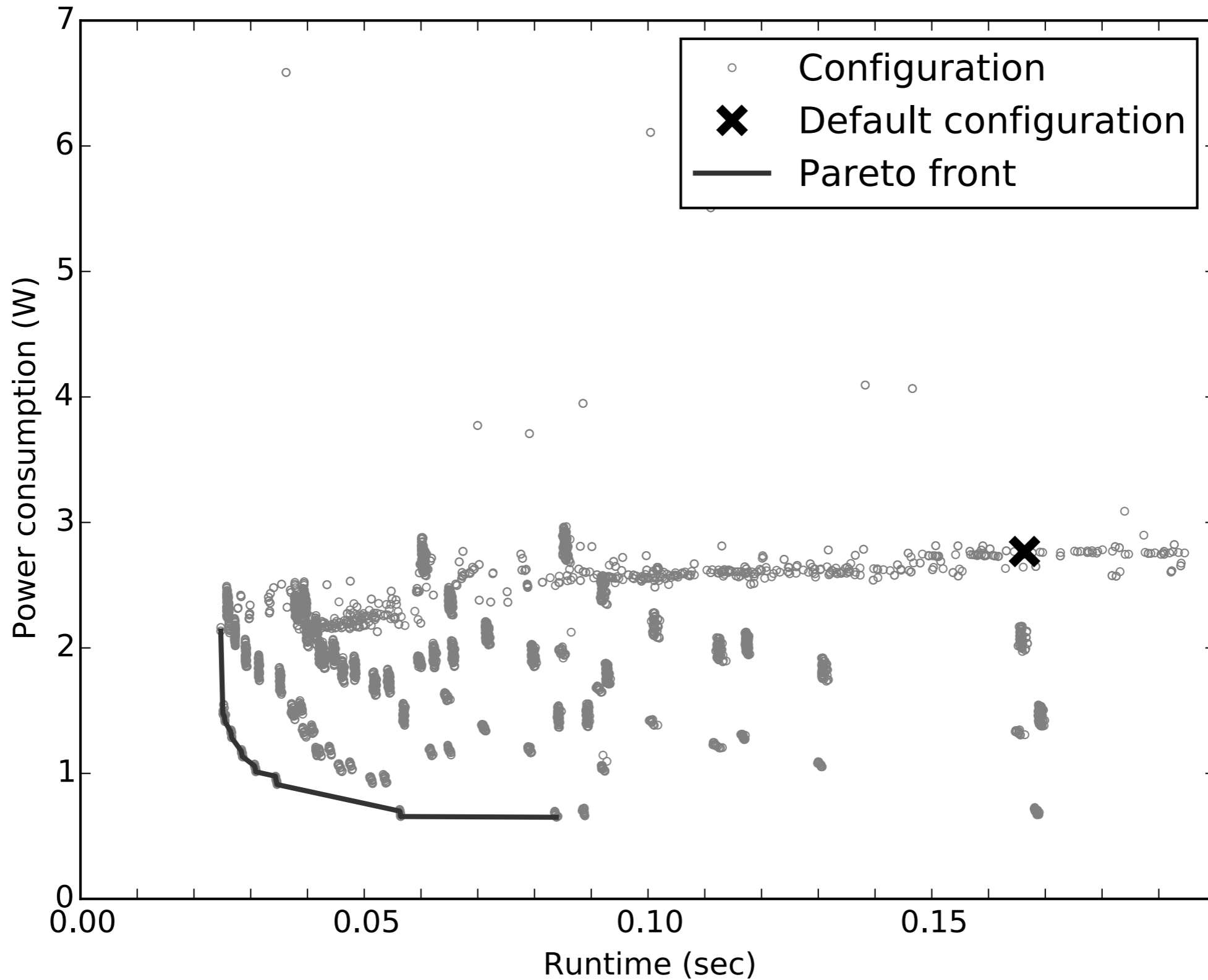
Black square: a parameter which when increased worsen the metric



Results KinectFusion compiler DSE speedup



Results KinectFusion architecture DSE power/runtime



DSE final result

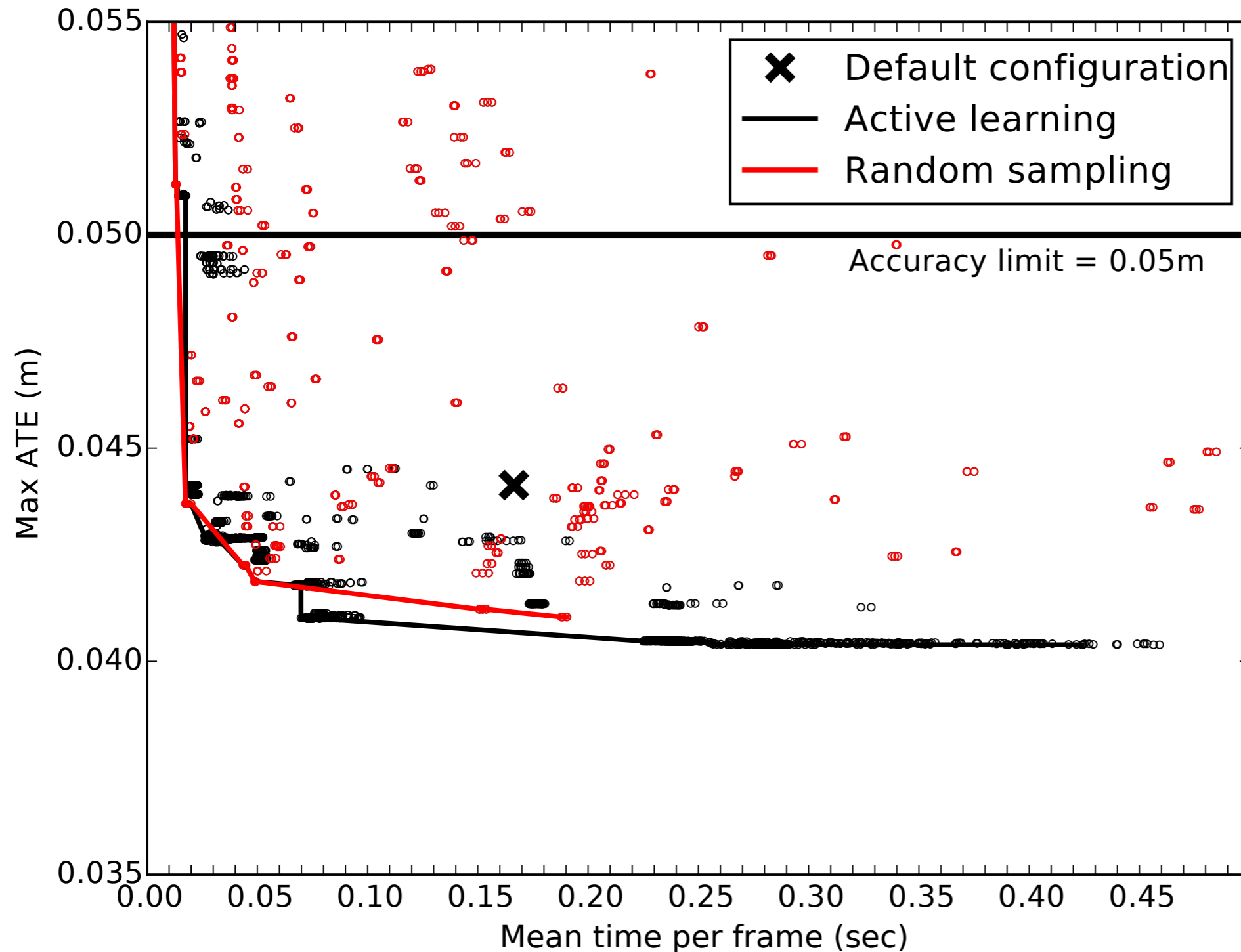
Constraint	Runtime (FPS)	Max ATE (cm)	Power (Watts)
Default	6.03	4.41	2.77
Best runtime	39.85	4.47	1.47
Best accuracy	1.51	3.30	2.38
Best power	11.92	4.45	0.65
Power < 1W	29.09	4.47	0.98
Power < 2W	39.85	4.47	1.47
FPS > 10	11.92	4.45	0.65
FPS > 20	28.87	4.47	0.91
FPS > 30	32.38	4.47	1.01

- Most of the improvement comes from the algorithmic space
- KinectFusion real-time on a popular embedded device
- Enabling auto-tuning at the domain-specific language (DSL) level



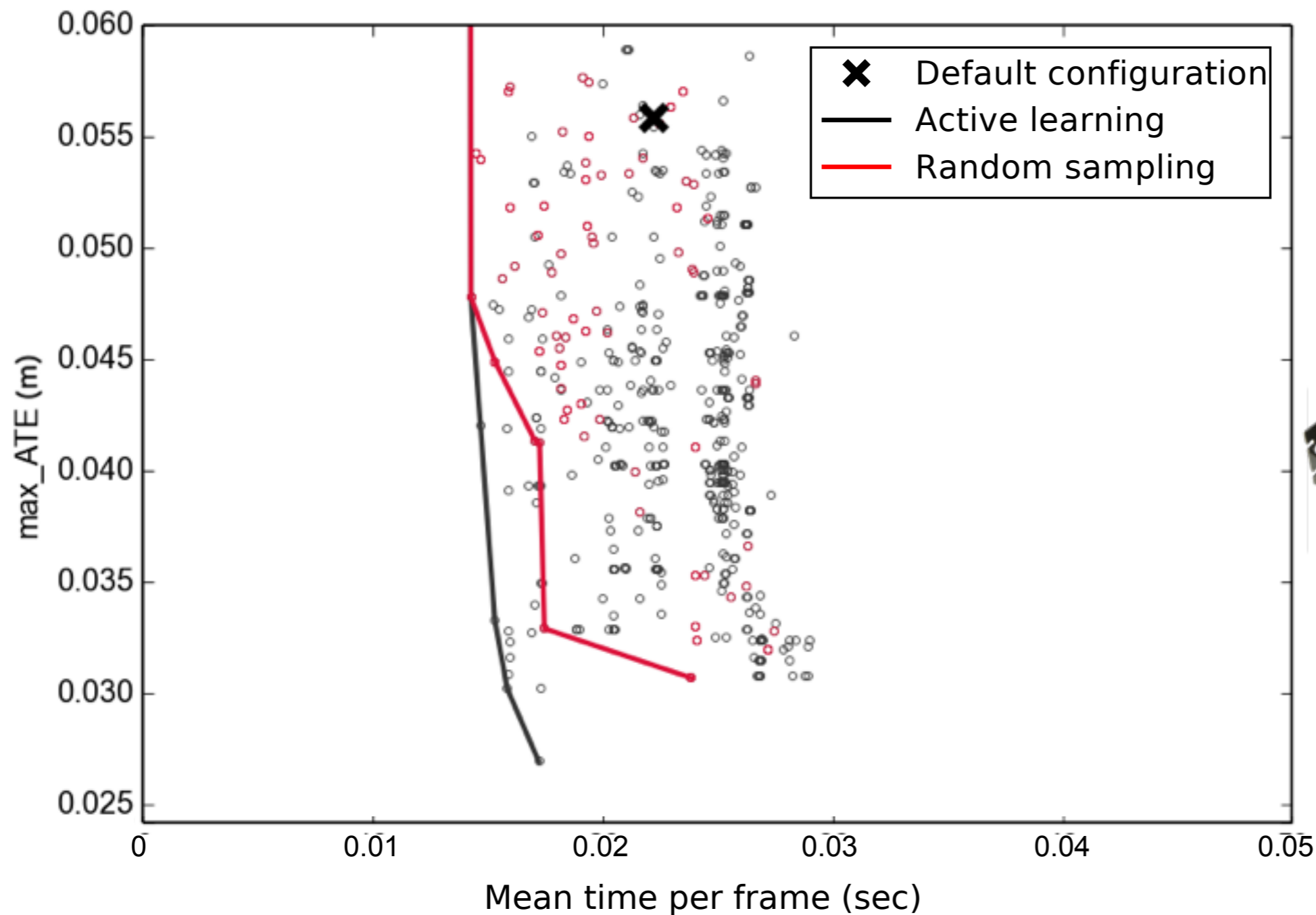
Results KinectFusion algo DSE error/runtime

Machine	Type	CPU	CPU name	CPU cores	GPU	GPU name
ASUS T200TA	Detachable laptop	Intel Silvermont	Atom Z3795	4	Intel	HD Graphics



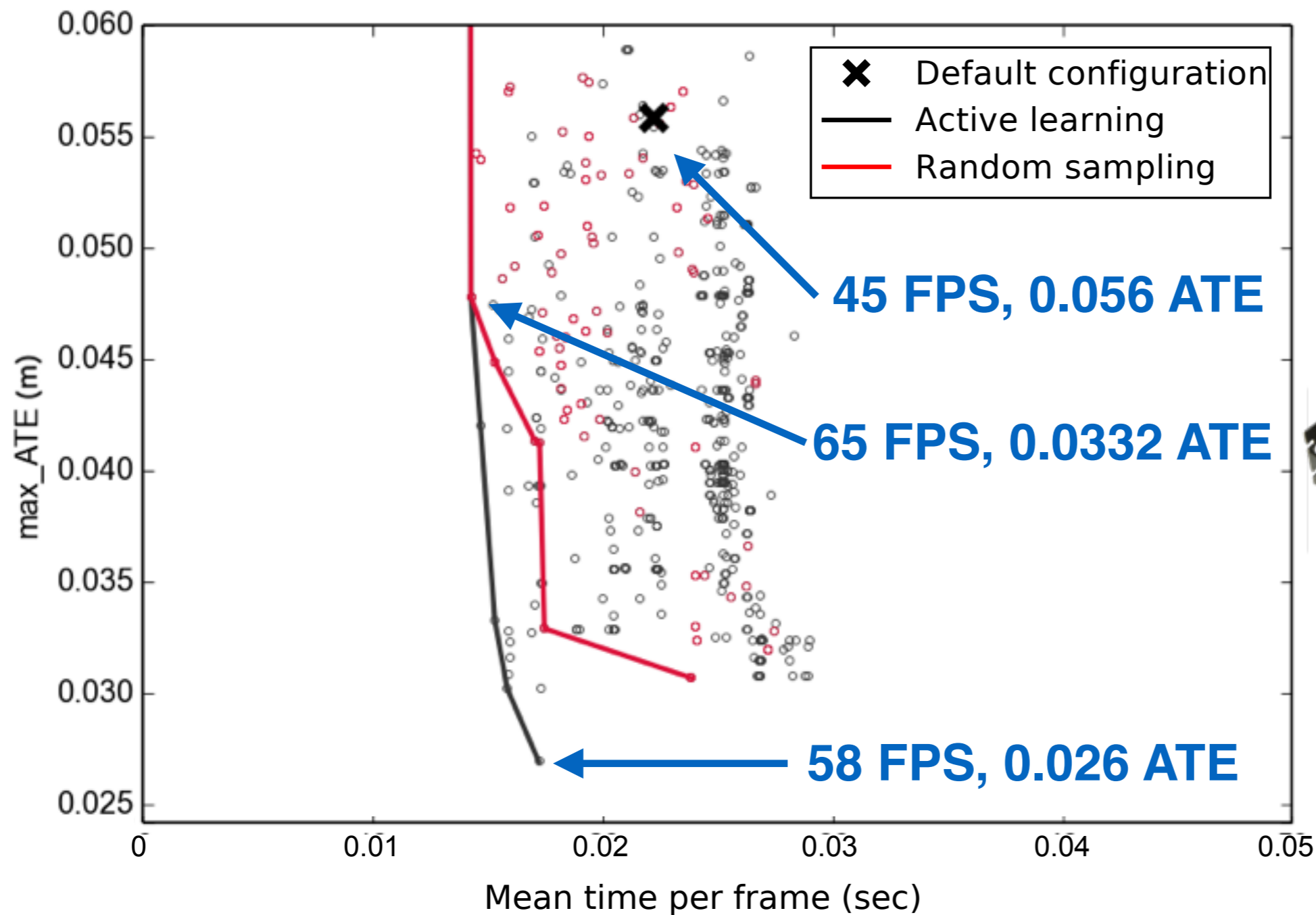
Results **ElasticFusion** algo DSE error/runtime

Machine	Type	CPU	CPU name	CPU cores	GPU	GPU name
NVIDIA/Intel	Desktop	Intel Ivy Bridge	E5-1620	8	NVIDIA	GTX 780 Ti

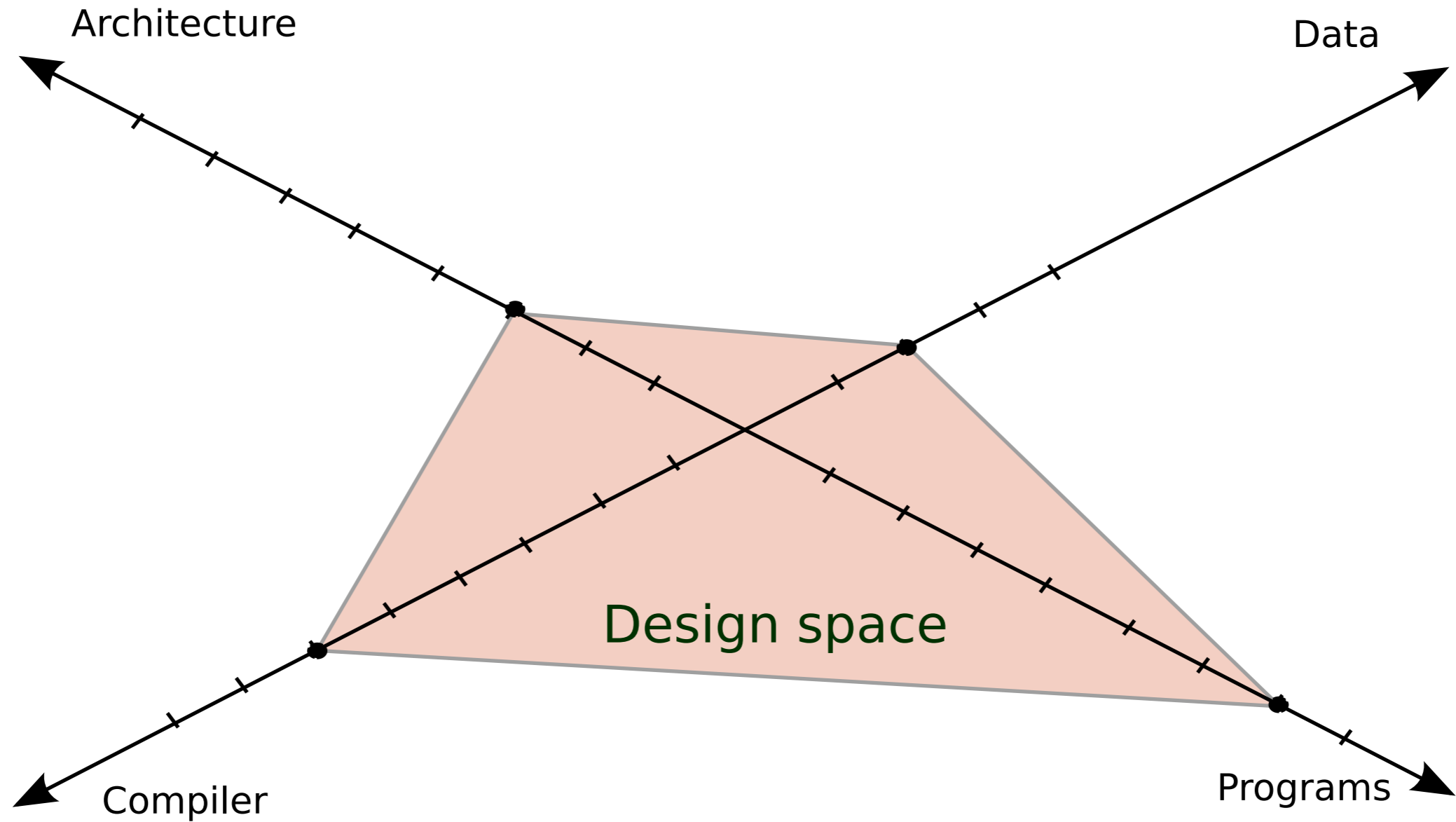


Results **ElasticFusion** algo DSE error/runtime

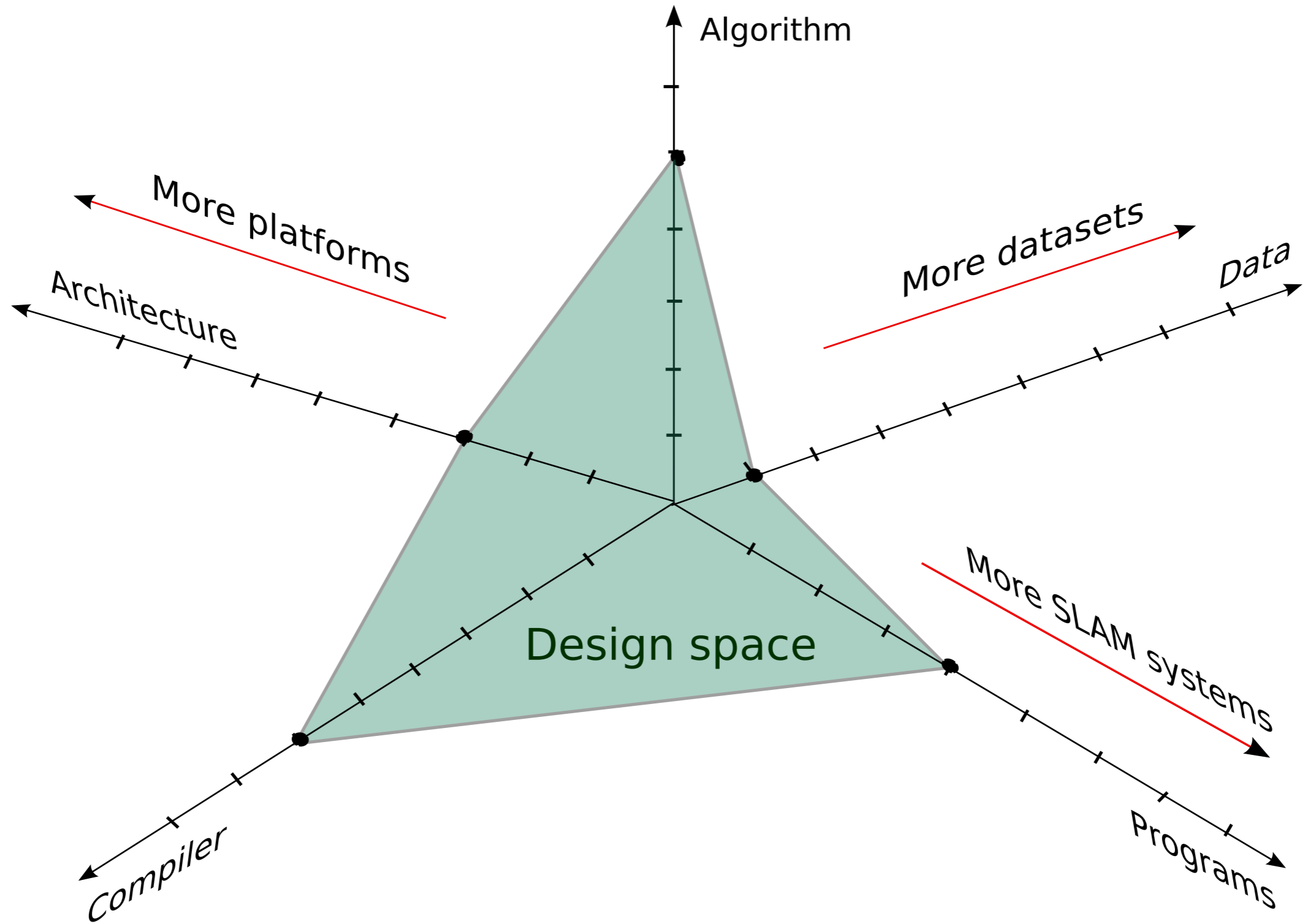
Machine	Type	CPU	CPU name	CPU cores	GPU	GPU name
NVIDIA/Intel	Desktop	Intel Ivy Bridge	E5-1620	8	NVIDIA	GTX 780 Ti



DSE the big picture I



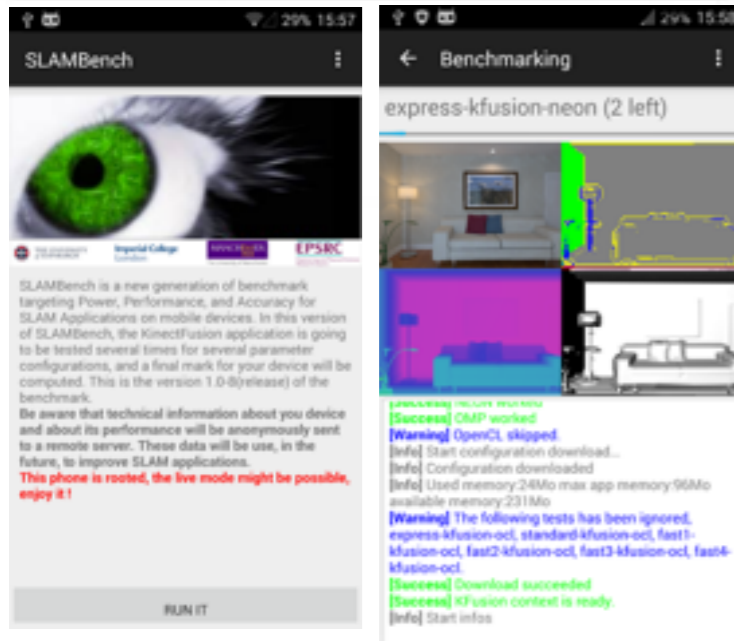
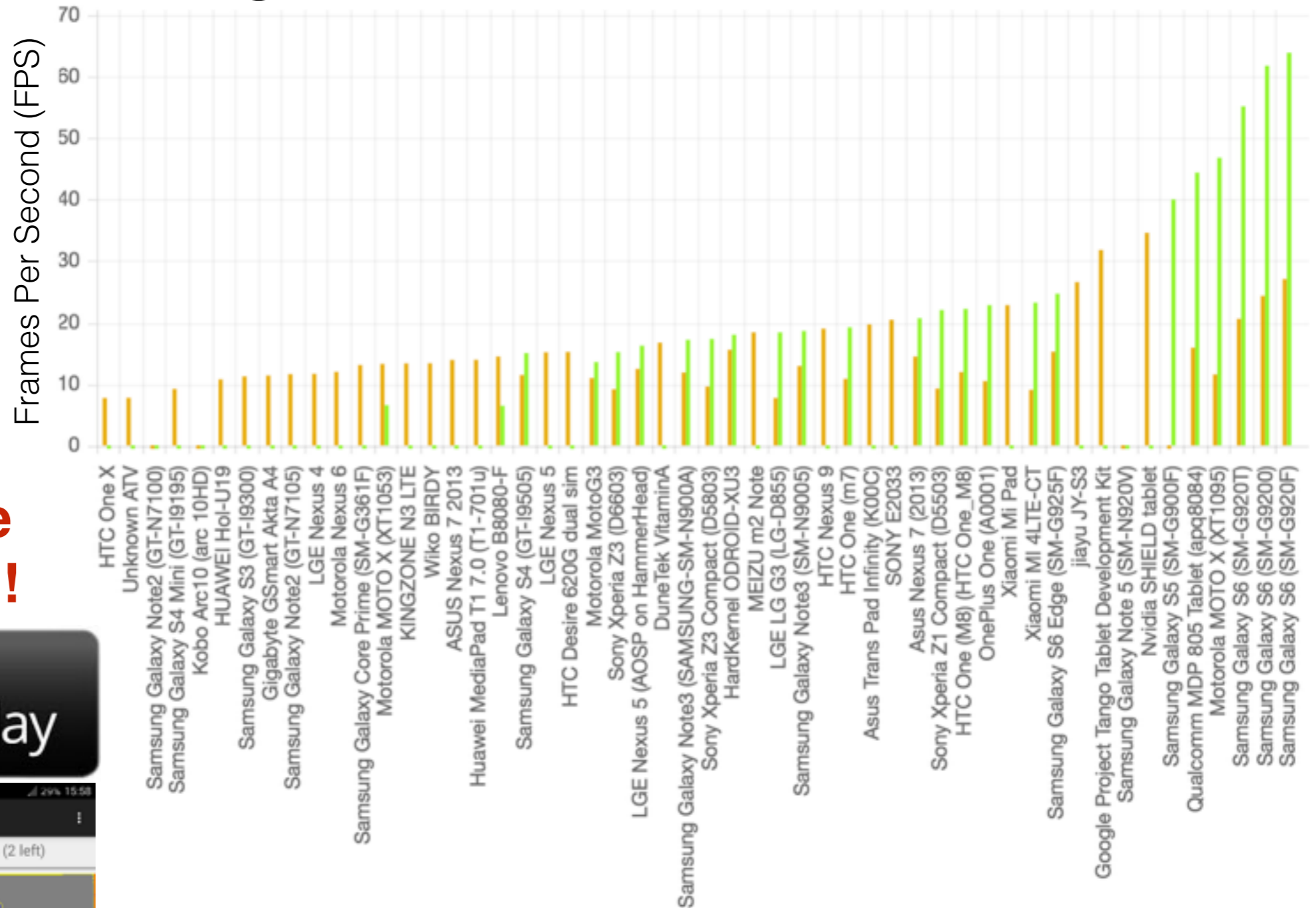
DSE the big picture II



Crowdsourcing mobile Android SLAMBench

- SLAMBench OpenMP
- SLAMBench OpenCL

Get it now,
And see where
your device is!!

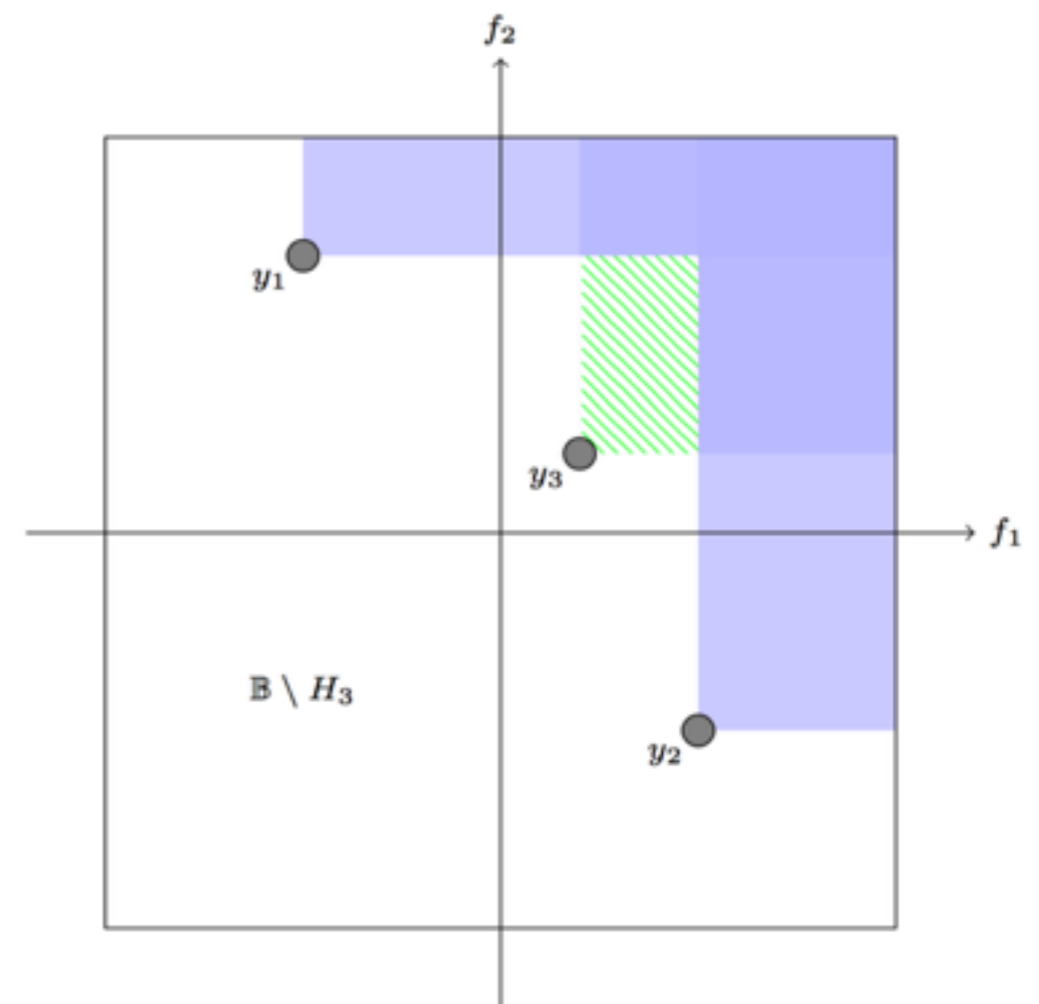
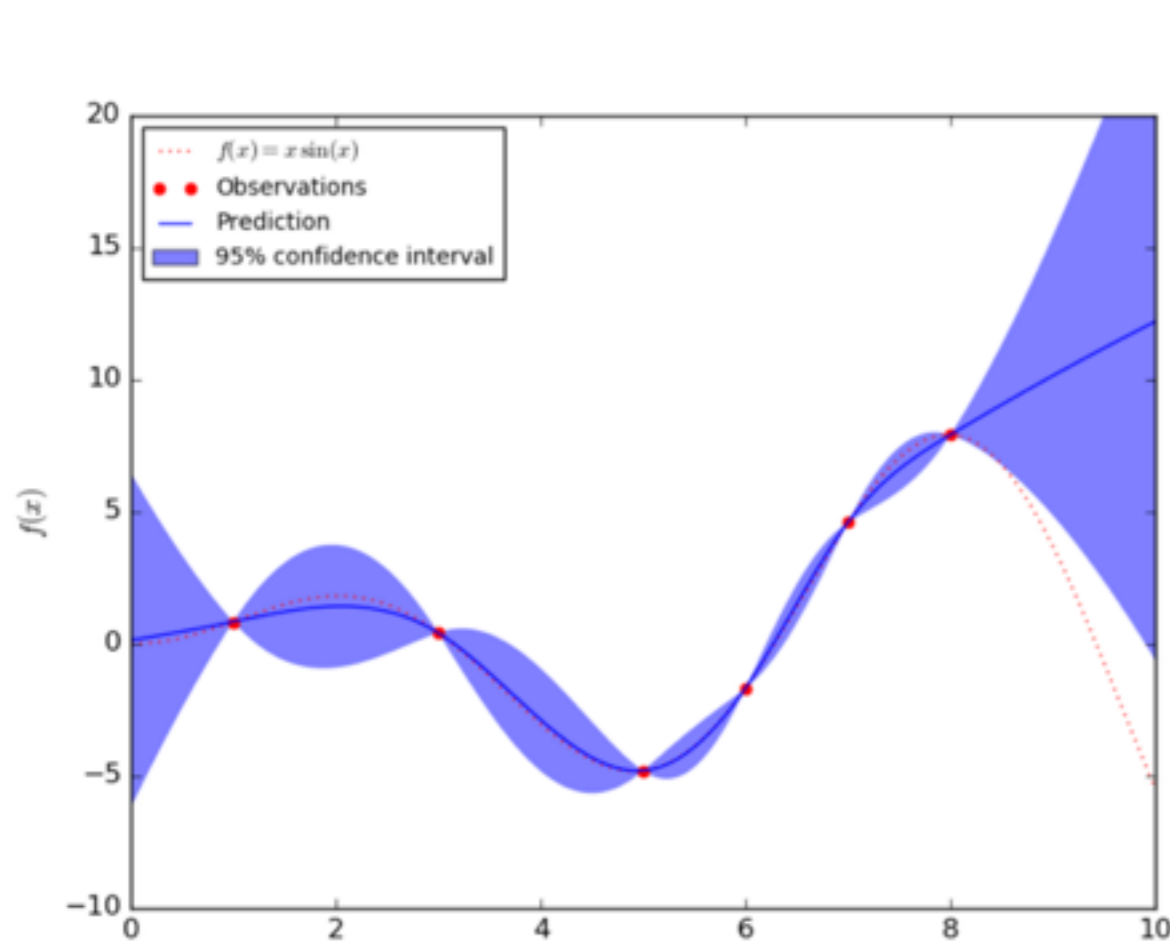


- It runs a set of configurations on the available languages on your device
- Then shows the best achieved result



Future work

- Gaussian Processes (GPs)
- Global optimisation of Black-Box functions
- Multi-objective optimisation: Expected HyperVolume Improvement (EHVI)



A Bayesian approach to constrained single- and multi-objective optimization, Paul FELIOT, Julien BECT and Emmanuel VAZQUEZ, 2015

Conclusion - take away messages

1. Building tools to explore the performance landscape for SLAM solutions
2. Generalisation to other applications
3. Multi-objective optimisation: speed/power/accuracy
4. Semantic accuracy check is very powerful:
 - enables non bit-wise accuracy check
 - aggressive approximate computing and auto-tuning
5. Pareto maps how configurations should be adapted when objectives change - static and dynamic
6. Large improvement over default configuration



References I

- [Nardi et al. 2015] L. Nardi, B. Bodin, M. Z. Zia, J. Mawer, A. Nisbet, P. H. J. Kelly, A. J. Davison, M. Luján, M. F. P. O'Boyle, G. Riley, N. Topham, and S. Furber. "Introducing SLAMBench, a performance and accuracy benchmarking methodology for SLAM." Submitted, arXiv:1410.2167, 2015.
- [Newcombe et al. ICCV 2011] R. A. Newcombe, S. J. Lovegrove and A. J. Davison. "DTAM: Dense tracking and mapping in real-time." Computer Vision (ICCV), 2011 IEEE International Conference on. IEEE, 2011.
- [Rusinkiewicz and Levoy 2001] S. Rusinkiewicz, and M. Levoy. "Efficient variants of the ICP algorithm." 3-D Digital Imaging and Modeling, 2001. Proceedings. Third International Conference on. IEEE, 2001.
- [Chen et al. 2013] J. Chen, D. Bautembach, and S. Izadi, Scalable real-time volumetric surface reconstruction, in ACM Trans. Graph., 2013.
- [Newcombe et al. ISMAR 2011] R. A. Newcombe, S. Izadi, O. Hilliges, D. Molyneaux, D. Kim, A. J. Davison, P. Kohi, J. Shotton, S. Hodges, and A. Fitzgibbon. "KinectFusion: Real-time dense surface mapping and tracking." 10th IEEE Int. Symp. on Mixed and augmented reality (ISMAR), 2011.
- [Handa et al. 2014] A. Handa, T. Whelan, J. McDonald, and A. J. Davison. A Benchmark for RGB-D Visual Odometry, 3D Reconstruction and SLAM. IEEE Int. Conf. on Robotics and Automation, ICRA 2014.
- [Reitmayr] G. Reitmayr. KFusion github 2011. <https://github.com/GerhardR/kfusion>
- [Curless and Levoy 1996] B. Curless and M. Levoy. A volumetric method for building complex models from range images. In Proc. Computer graphics and interactive technique. ACM, 1996.
- [Whelan et al. 2012] T. Whelan, M. Kaess, M. Fallon, H. Johannsson, J. Leonard, and J. McDonald. Kintinuous: Spatially extended kinectfusion. 2012.
- C. Jiawen, D. Bautembach, and S. Izadi. "Scalable real-time volumetric surface reconstruction." ACM TOG, 2013.
- Frahm, Jan-Michael, et al. "Building Rome on a cloudless day." Computer Vision–ECCV 2010. Springer Berlin Heidelberg, 2010.
- Erhan, Dumitru, et al. "Scalable object detection using deep neural networks." Proceedings of the IEEE CVPR. 2014.



References II

- Arbelaez, Pablo, et al. "Contour detection and hierarchical image segmentation." IEEE Pattern Analysis and Machine Intelligence, 2011.
- [Ogilvie 2014] Ogilvie, William, et al. "Fast automatic heuristic construction using active learning." Proceedings of the Workshop on Languages and Compilers for Parallel Computing (LCPC'14). 2014.
- [Siegmund 2015] Siegmund Norbert et al. "Performance-influence models for highly configurable systems", submitted FSE 2015.
- [Guo 2013] Guo, Jianmei, et al. "Variability-aware performance prediction: A statistical learning approach." Automated Software Engineering (ASE), 2013 IEEE/ACM 28th International Conference on. IEEE, 2013.
- [Grewe 2011] Grewe, Dominik et al. "A static task partitioning approach for heterogeneous systems using OpenCL." Compiler Construction. Springer Berlin Heidelberg, 2011.
- [Kurek 2013] Kurek, Maciej, Tianchi Liu, and Wayne Luk. "MULTI-OBJECTIVE SELF-OPTIMIZATION OF RECONFIGURABLE DESIGNS WITH MACHINE LEARNING." 2nd Workshop on Self-Awareness in Reconfigurable Computing Systems (SRCS'13). 2013.
- [Balaprakash 2013] Balaprakash, Prasanna, Robert B. Gramacy, and Stefan M. Wild. "Active-learning-based surrogate models for empirical performance tuning." Cluster Computing (CLUSTER), 2013 IEEE International Conference on. IEEE, 2013.
- [Vespa 2015] Vespa Emanuele. "Sparse voxelization of dense volumetric reconstruction with automated analysis of scene reconstruction quality." M.Res. thesis, Imperial College London, 2015.

