

1st TVM and Deep Learning Compilation Conference

tvm

AMPL

December 12, 2018

PAUL G.
ALLEN
SCHOOL

W

Luis Ceze

Welcome to the 1st TVM and Deep Learning Compilation Conference!

Welcome to the 1st TVM and Deep Learning Compilation Conference!

180+ ppl!

Machine learning is amazing...

Machine learning is amazing...

super human accuracy,
self driving cars, automated scientific
discoveries...

Machine learning is amazing...

super human accuracy,
self driving cars, automated scientific
discoveries...

wow!

Software era:

Problem to solve

Write code

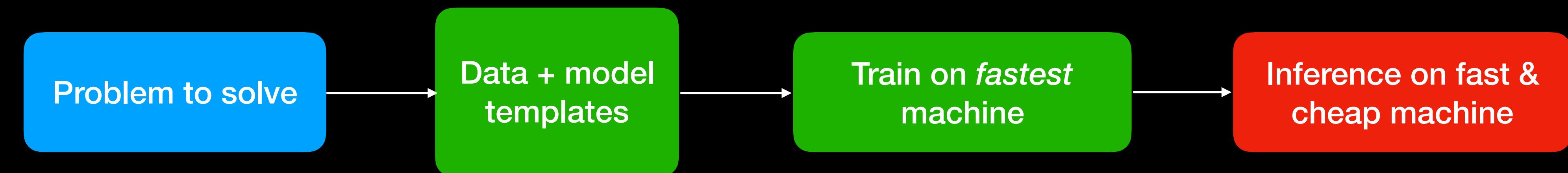
Run on fast machine



Software era:



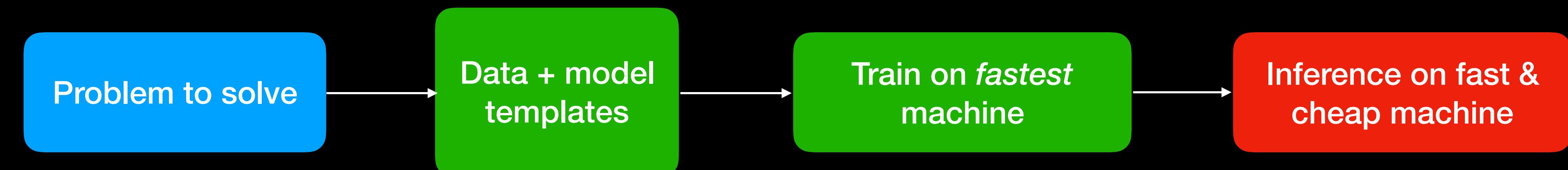
Machine learning era:



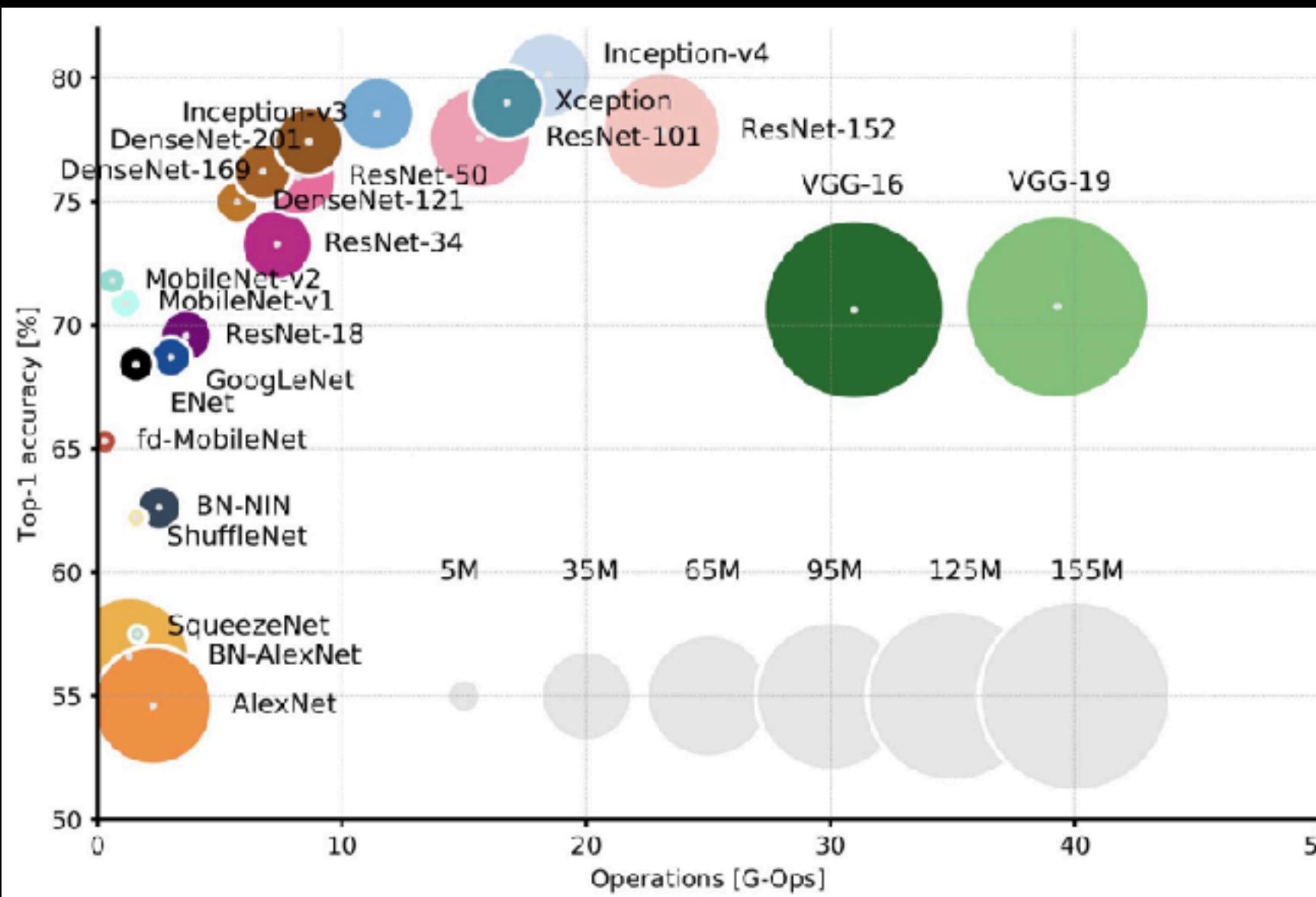
Software era:



Machine learning era:



Model size and compute cost growing fast



Software era:

Problem to solve

Write code

Run on fast machine

Machine learning era:

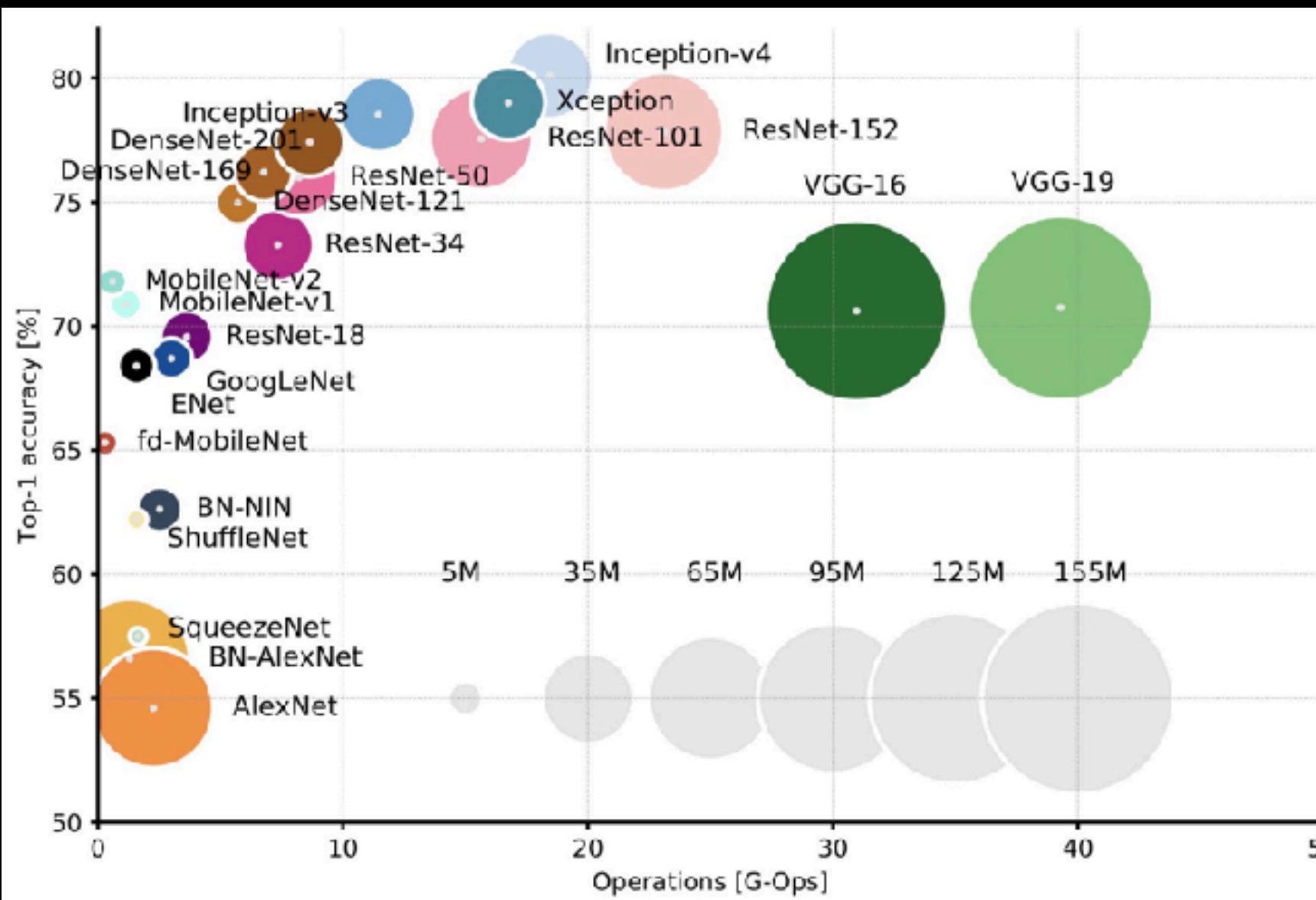
Problem to solve

Data + model templates

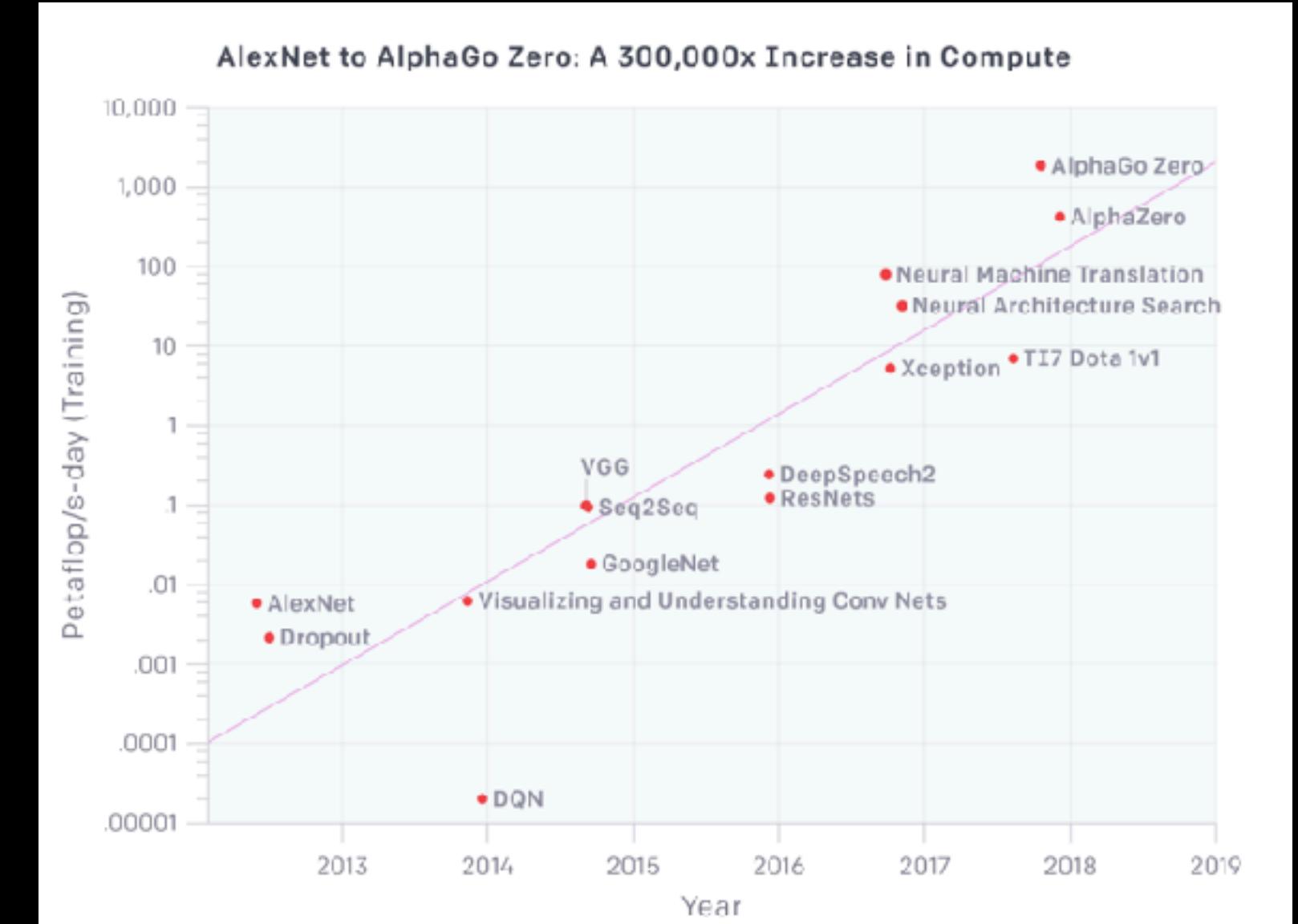
Train on *fastest* machine

Inference on fast & cheap machine

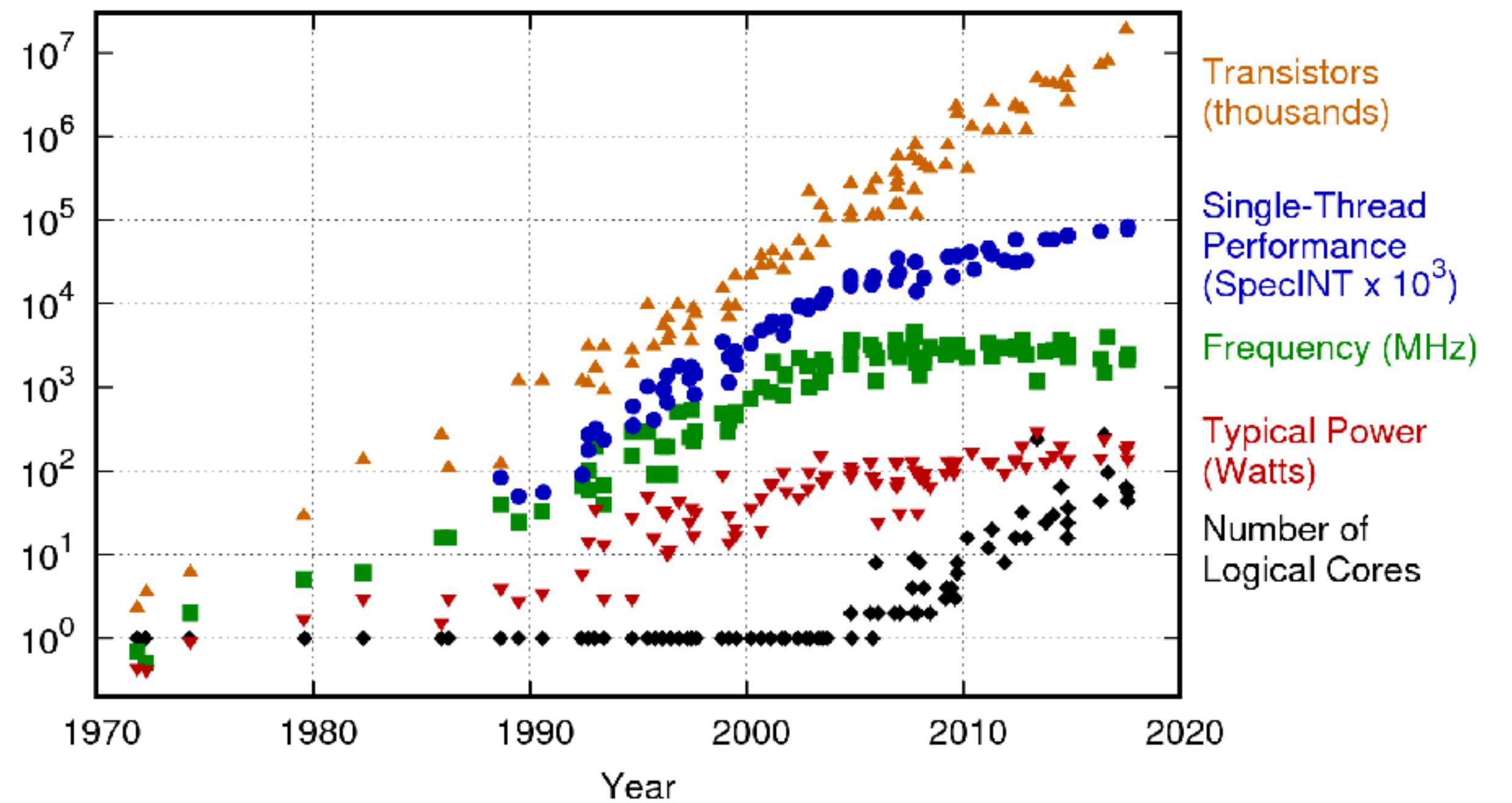
Model size and compute cost growing fast



Training costs growing exponentially

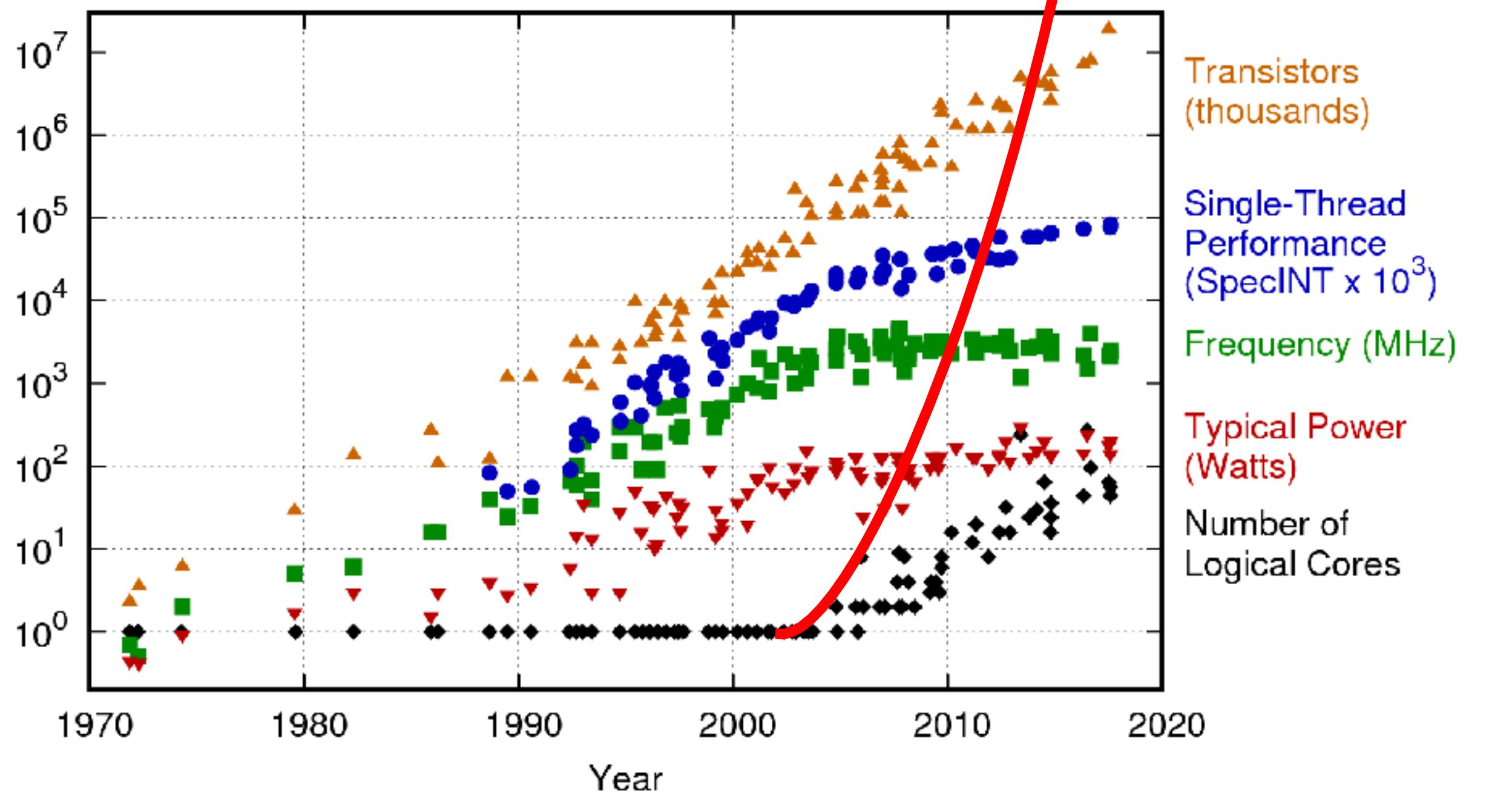


42 Years of Microprocessor Trend Data



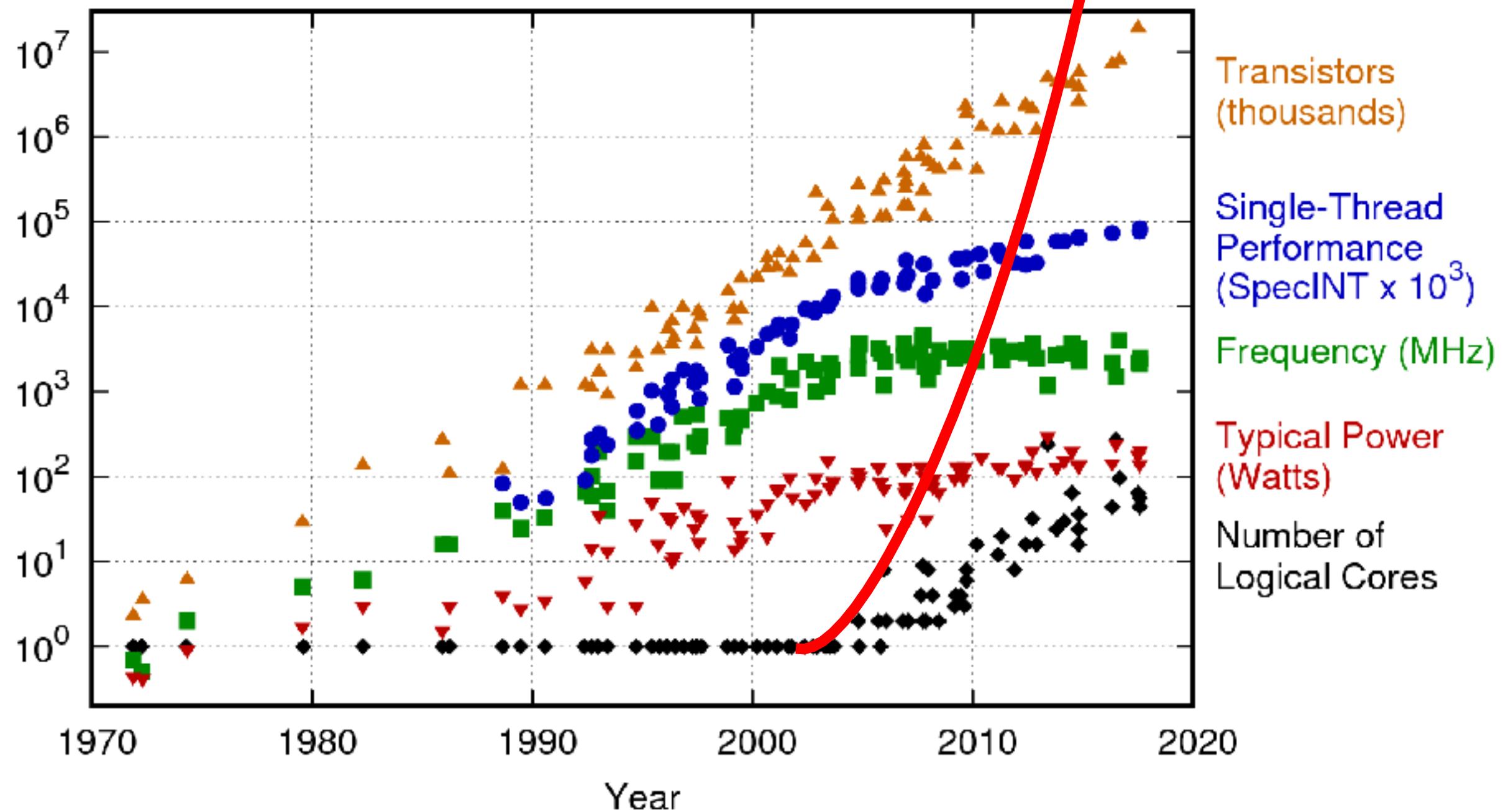
Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten
New plot and data collected for 2010-2017 by K. Rupp

42 Years of Microprocessor Trend Data



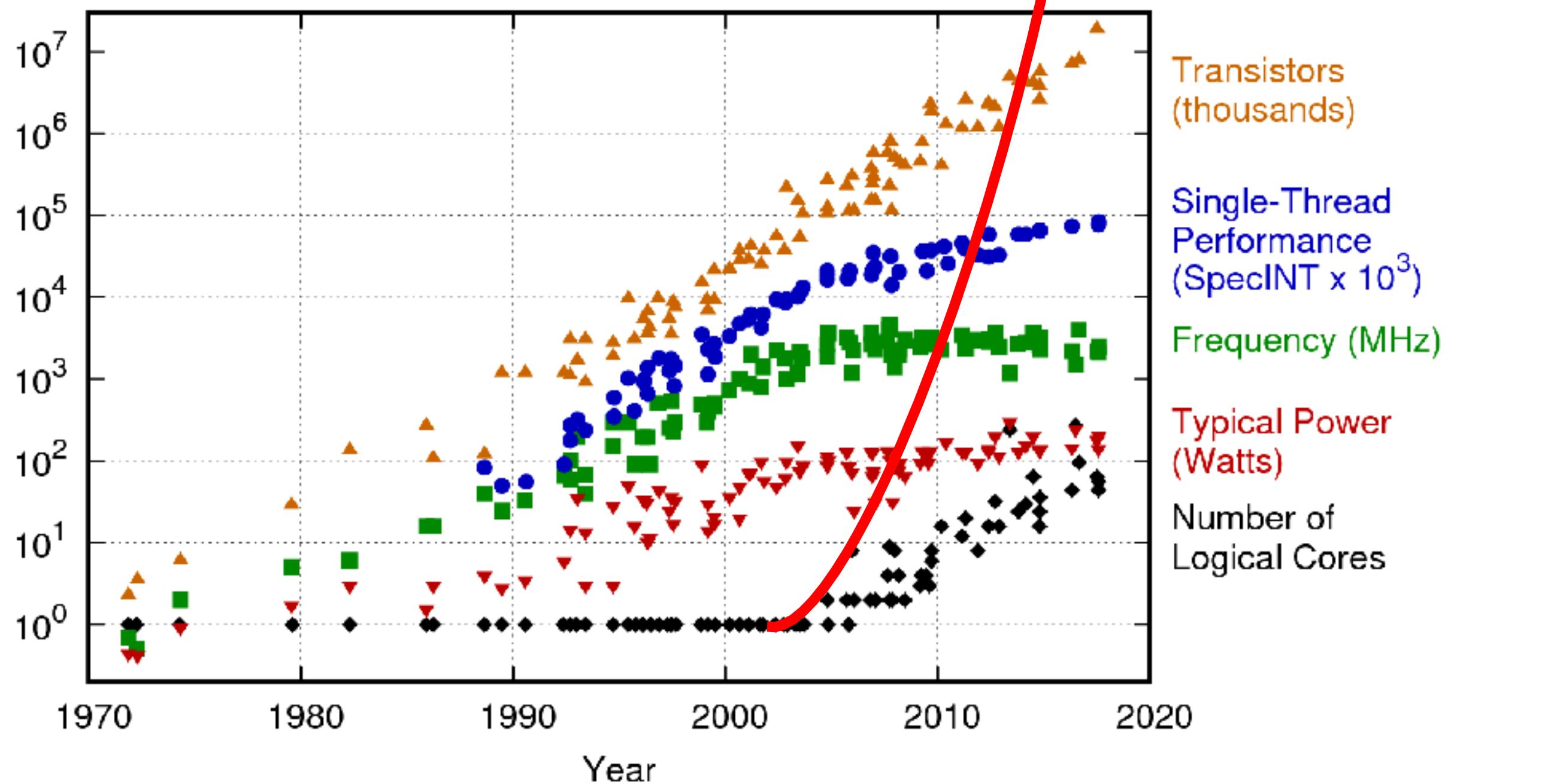
Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten
New plot and data collected for 2010-2017 by K. Rupp

42 Years of Microprocessor Trend Data

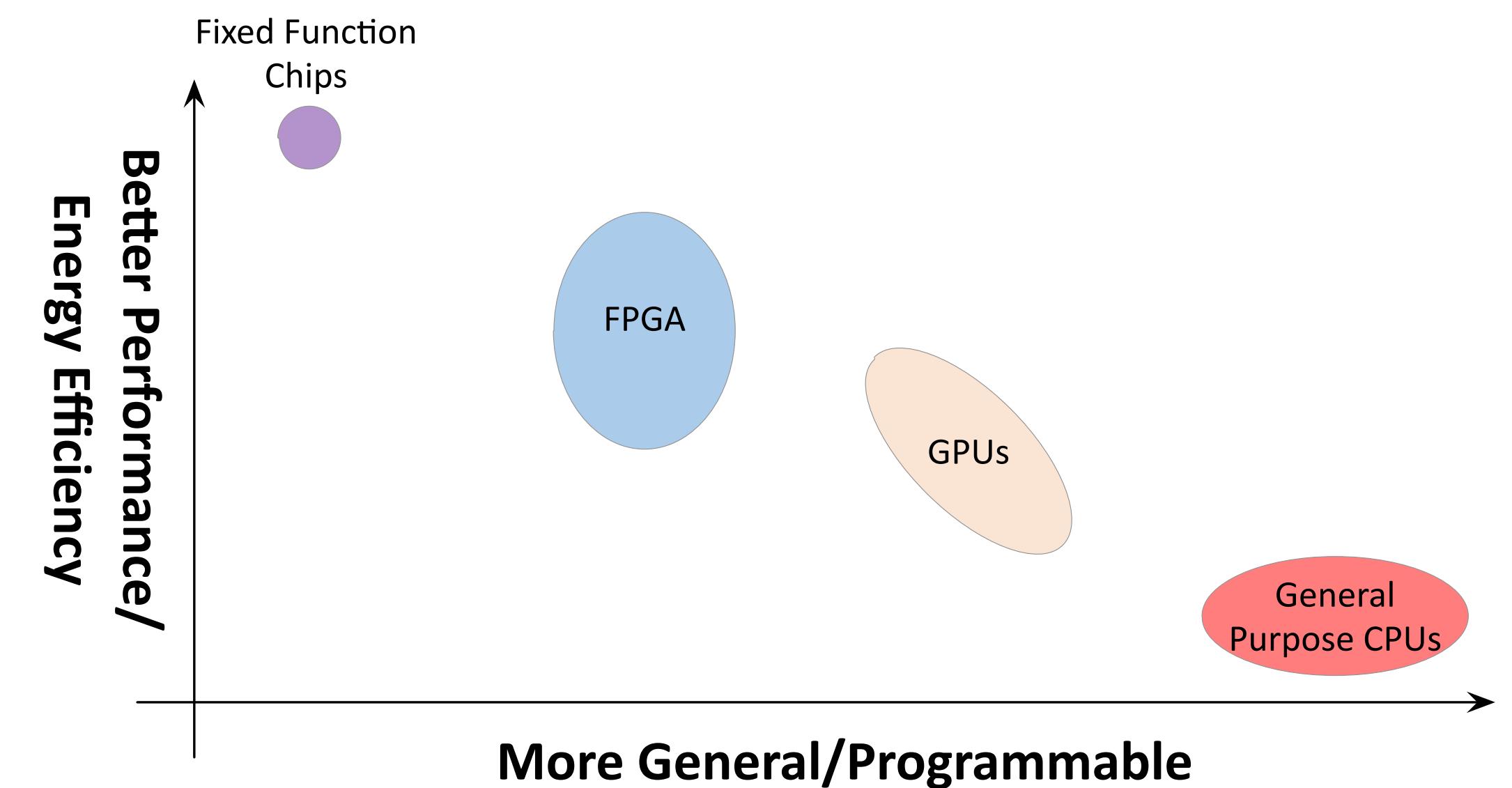


Fundamental trade-off between specialization and performance/efficiency.

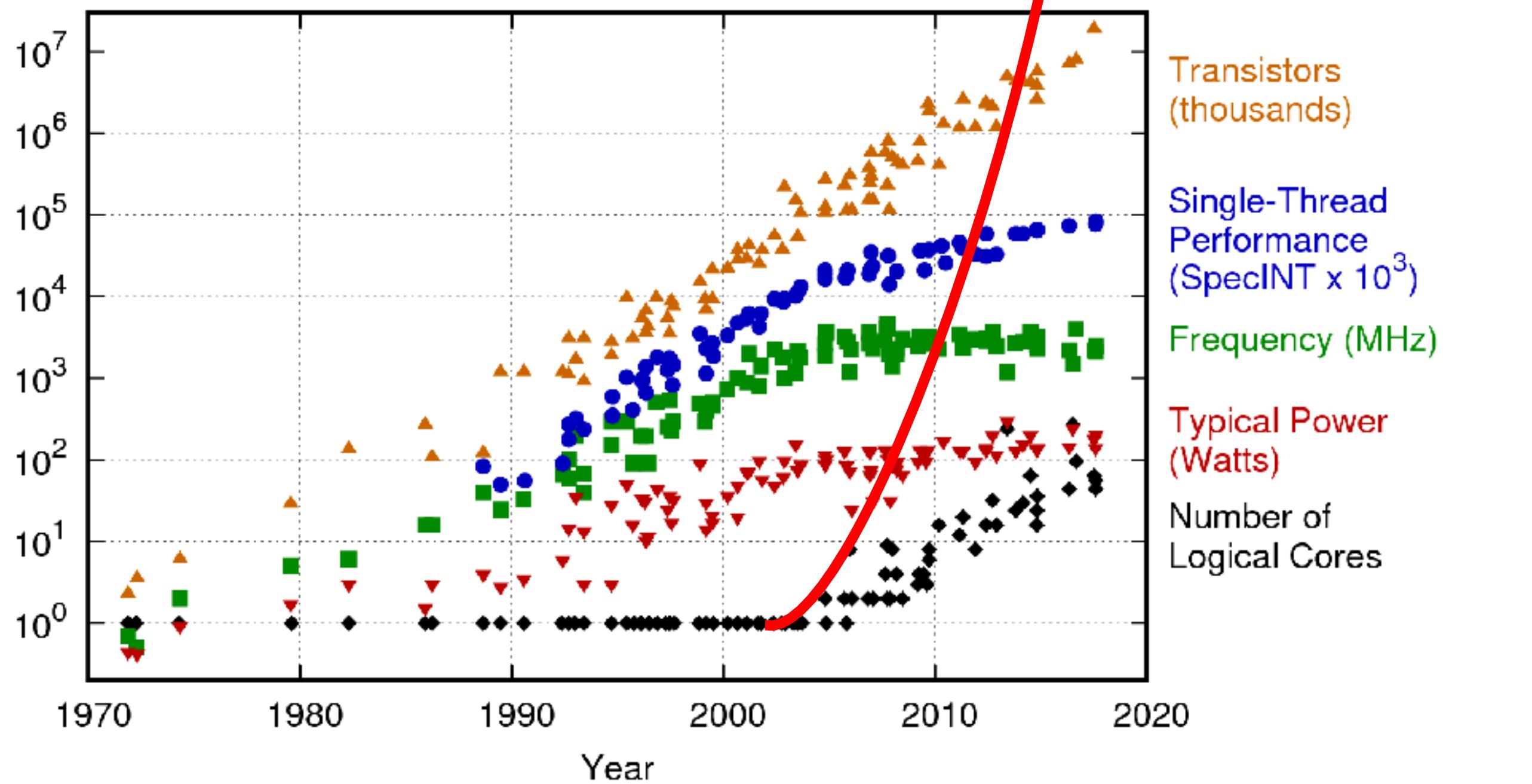
42 Years of Microprocessor Trend Data



Fundamental trade-off between specialization and performance/efficiency.

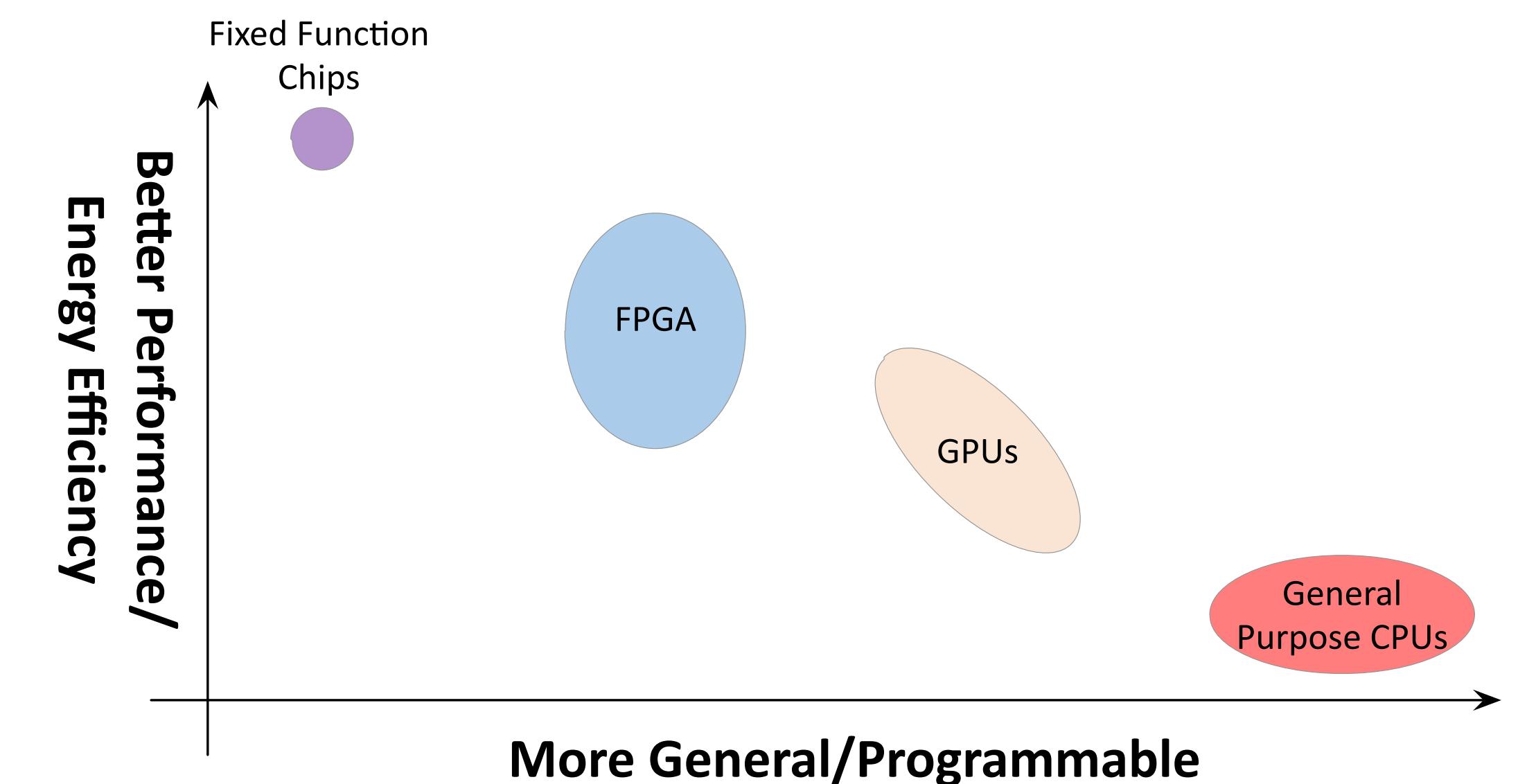


42 Years of Microprocessor Trend Data



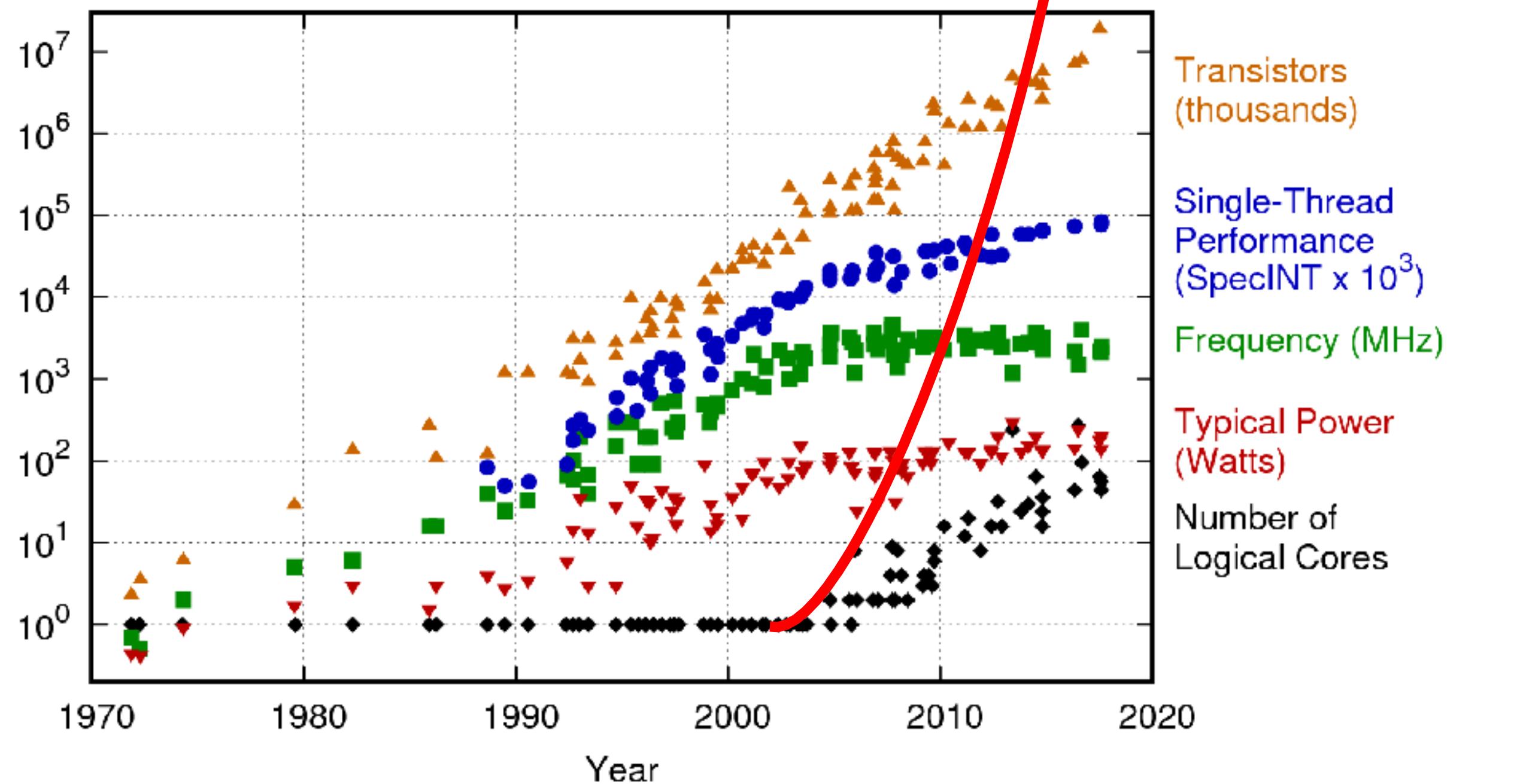
Popularity and computational cost of ML. Ops.

Fundamental trade-off between specialization and performance/efficiency.

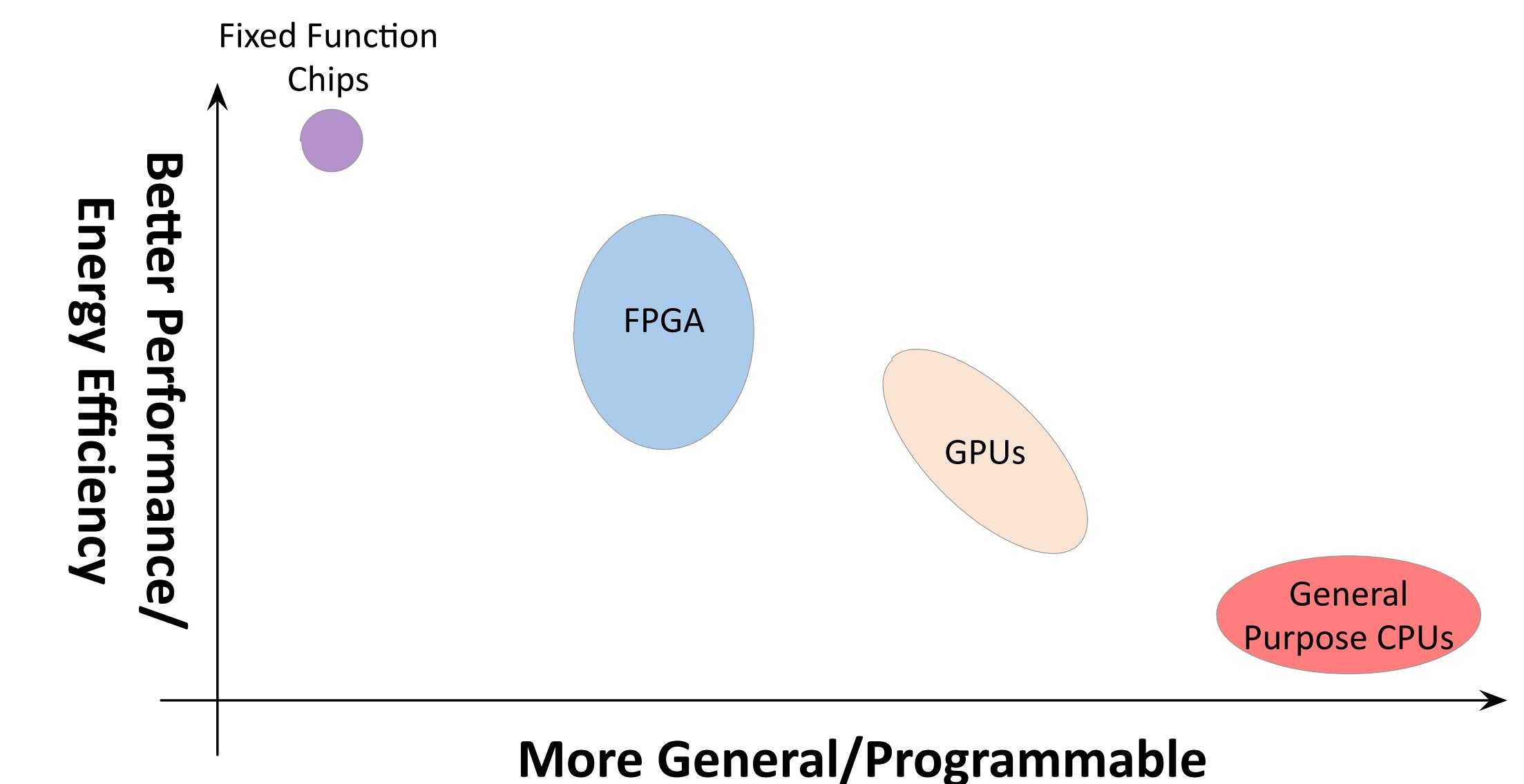


Machine learning algorithms are **relatively** simple to implement in HW... great!

42 Years of Microprocessor Trend Data



Fundamental trade-off between specialization and performance/efficiency.

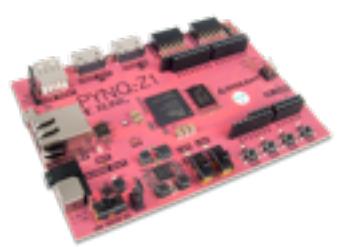


Machine learning algorithms are **relatively** simple to implement in HW... great!

Machine Learning Makes Computer Architecture Cool Again!

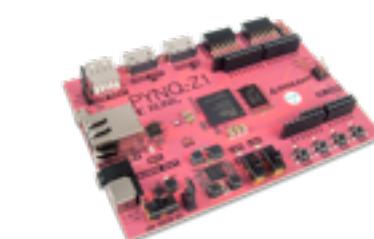


...





...



 XILINX

 Microsoft

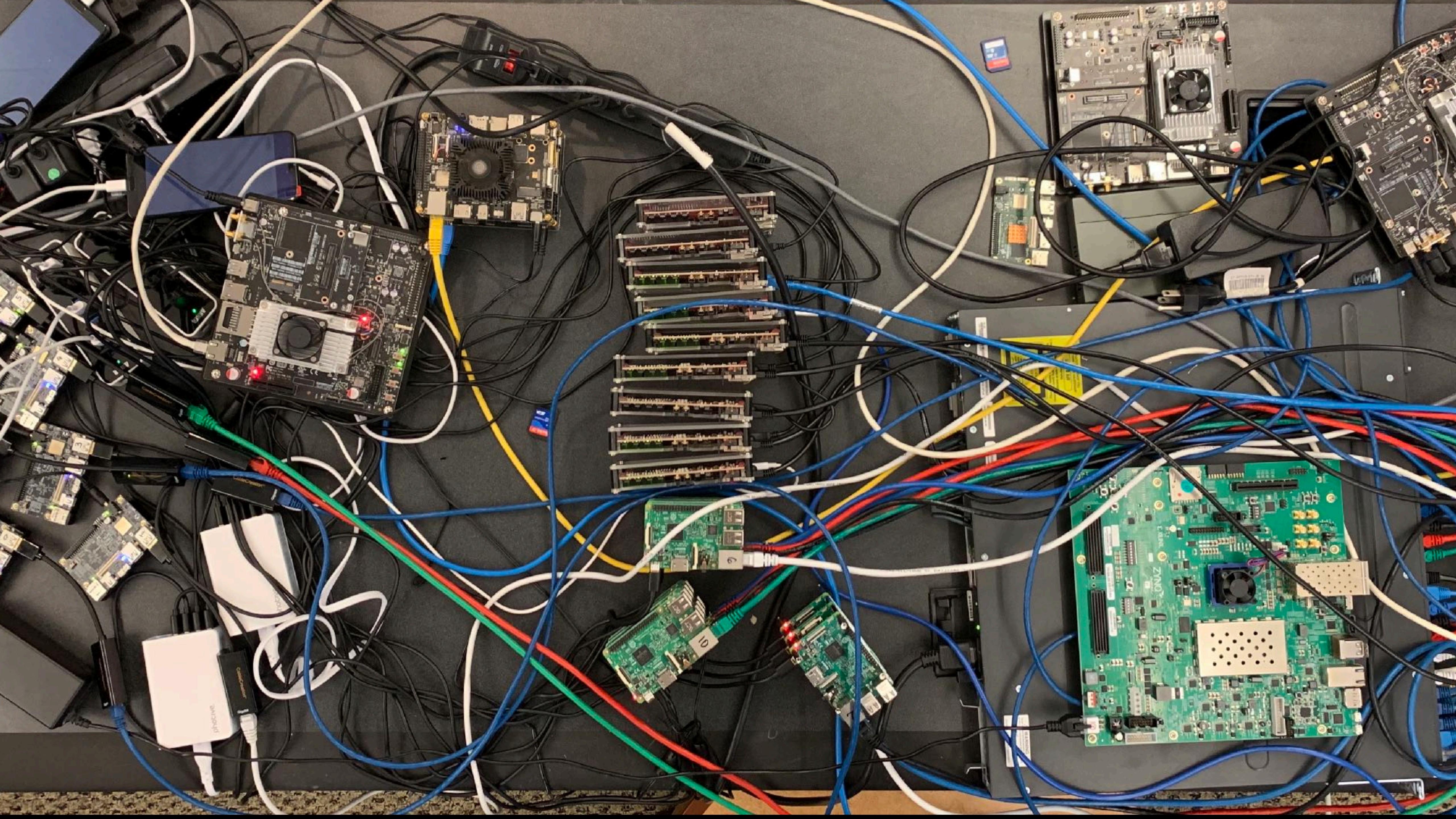
 QUALCOMM

 amazon

 Google

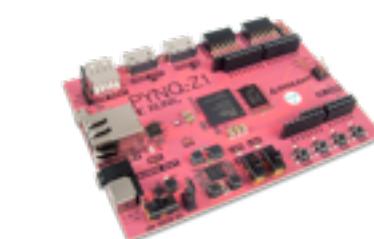
 HUAWEI

+~50 startups





...



 XILINX



+~50 startups

Models:

CNN

GAN

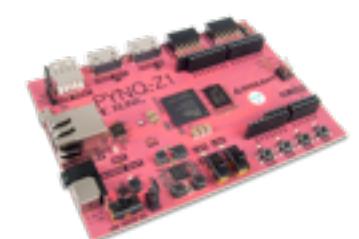
RNN

MLP

DQNN



...



+~50 startups

Models:

CNN

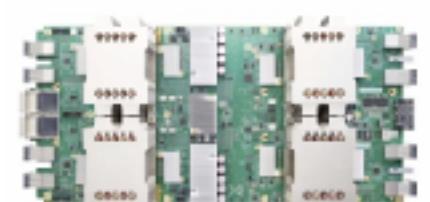
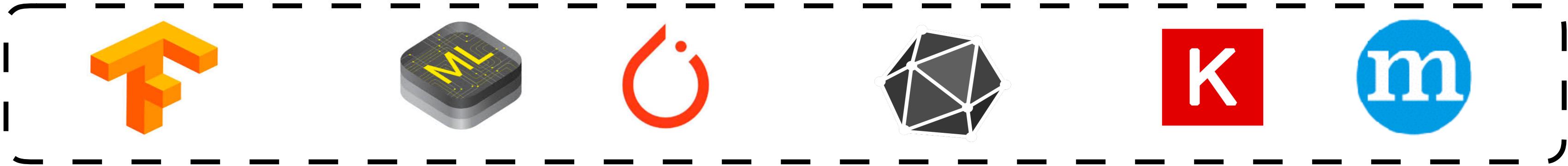
RNN

DQNN

GAN

MLP

Frameworks:



...



+~50 startups

Models:

CNN

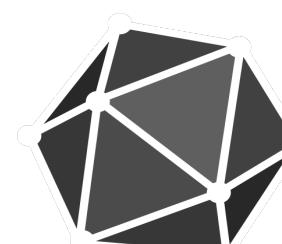
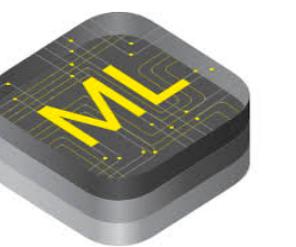
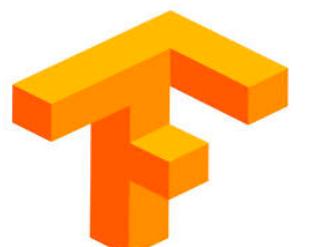
RNN

DQNN

GAN

MLP

Frameworks:



Challenge: Efficiently deploying deep learning everywhere



...



+~50 startups

Gaurav Kapoor, Core Machine Learning



HW+SW optimization is key for efficiency

HW+SW optimization is key for efficiency

Lots of hand-tuning, full automation would be a holy grail



Academic group focused on **S**ystems + Computer **A**rchitecture +
Machine **L**earning + **P**rogramming **L**anguages



Luis Ceze
Professor



Carlos Guestrin
Professor



Arvind Krishnamurthy
Professor



Zachary Tatlock
Assistant Professor



Meghan Cowan



Thierry Moreau



Eddie Yan



Luis Vega



Jared Roesch



Steven Lyubomirsky



Tianqi Chen



Haichen Shen



Liang Luo



Logan Weber



Pratyush Patel



Josh Fromm



Gus Smith



Ziheng Jiang



Marisa Kirisame



Lianmin Zheng



Seungyeop Han



Jacob Nelson



Amar Phanishayee



Academic group focused on **S**ystems + Computer **A**rchitecture +
Machine **L**earning + **P**rogramming **L**anguages



Luis Ceze
Professor



Carlos Guestrin
Professor



Arvind Krishnamurthy
Professor



Zachary Tatlock
Assistant Professor



Meghan Cowan



Thierry Moreau



Eddie Yan



Luis Vega



Jared Roesch



Steven Lyubomirsky



Tianqi Chen



Haichen Shen



Liang Luo



Logan Weber



Pratyush Patel



Josh Fromm



Gus Smith



Ziheng Jiang



Marisa Kirisame



Lianmin Zheng



Josh Pollock



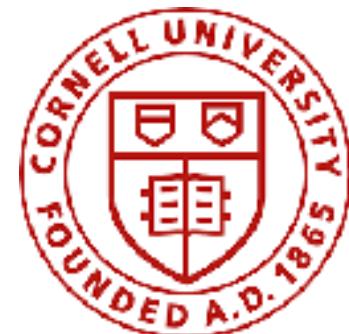
Seungyeop Han



Jacob Nelson



Amar Phanishayee



Ucla

Berkeley
UNIVERSITY OF CALIFORNIA

THE UNIVERSITY OF
TEXAS
AT AUSTIN



Academic group focused on **S**ystems + Computer **A**rchitecture +
Machine **L**earning + **P**rogramming **L**anguages



Luis Ceze
Professor



Carlos Guestrin
Professor



Arvind Krishnamurthy
Professor



Zachary Tatlock
Assistant Professor



Meghan Cowan



Thierry Moreau



Eddie Yan



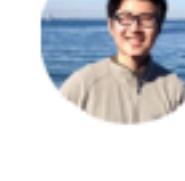
Luis Vega



Jared Roesch



Steven Lyubomirsky



Tianqi Chen



Haichen Shen



Liang Luo



Logan Weber



Pratyush Patel



Josh Fromm



Gus Smith



Ziheng Jiang



Josh Pollock



Marisa Kirisame



Lianmin Zheng



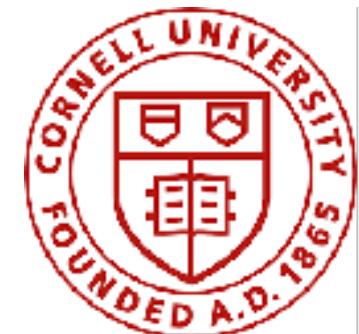
Seungyeop Han



Jacob Nelson



Amar Phanishayee



Ucla

Berkeley
UNIVERSITY OF CALIFORNIA

THE UNIVERSITY OF
TEXAS
AT AUSTIN



amazon
XILINX





PL: High-level support for future ML applications

Compilers: Extensible support for future models, optimizations and hardware architectures

Systems: On-device and cloud-based training, distributed systems for ML

Computer Architecture: Extensible, energy efficient hardware designs for inference and training



ML for Systems:
Automatic Learning-
Based Design and
Optimizations

PL: High-level support for future ML applications

Compilers: Extensible support for future models,
optimizations and hardware architectures

Systems: On-device and cloud-based training,
distributed systems for ML

Computer Architecture: Extensible, energy
efficient hardware designs for inference and training

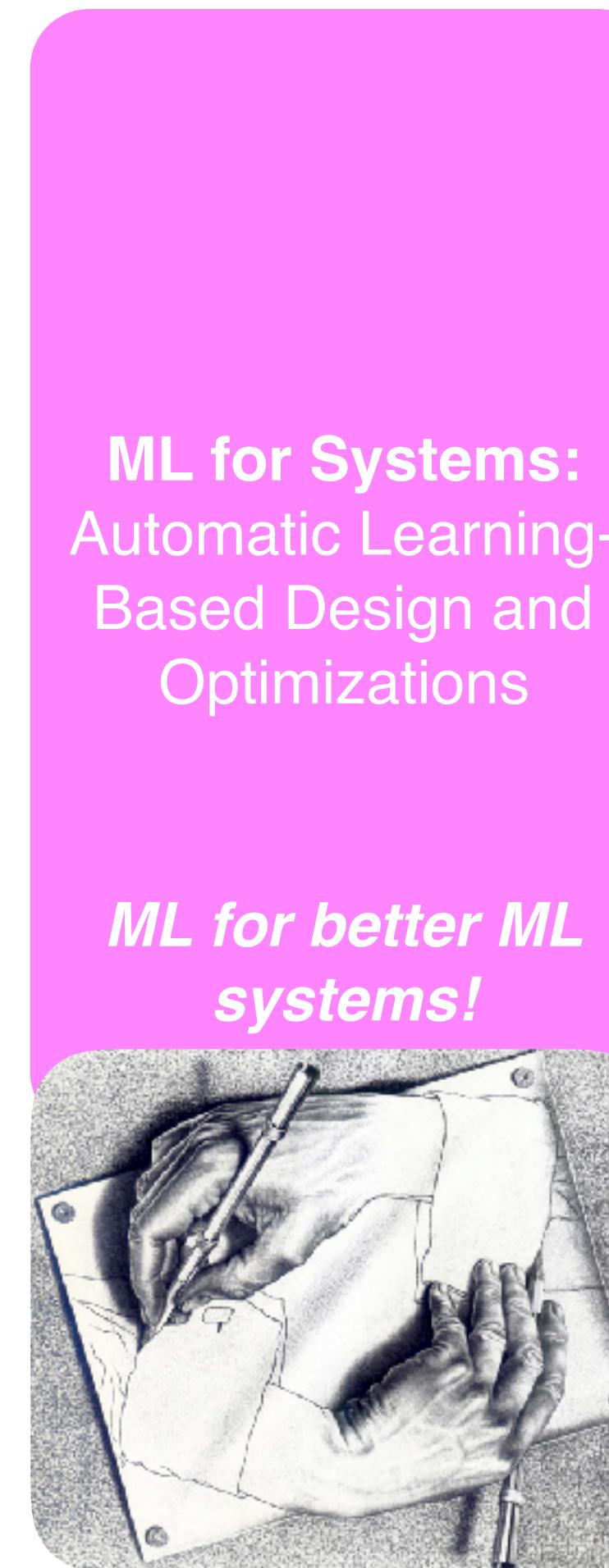


PL: High-level support for future ML applications

Compilers: Extensible support for future models, optimizations and hardware architectures

Systems: On-device and cloud-based training, distributed systems for ML

Computer Architecture: Extensible, energy efficient hardware designs for inference and training



PL: High-level support for future ML applications

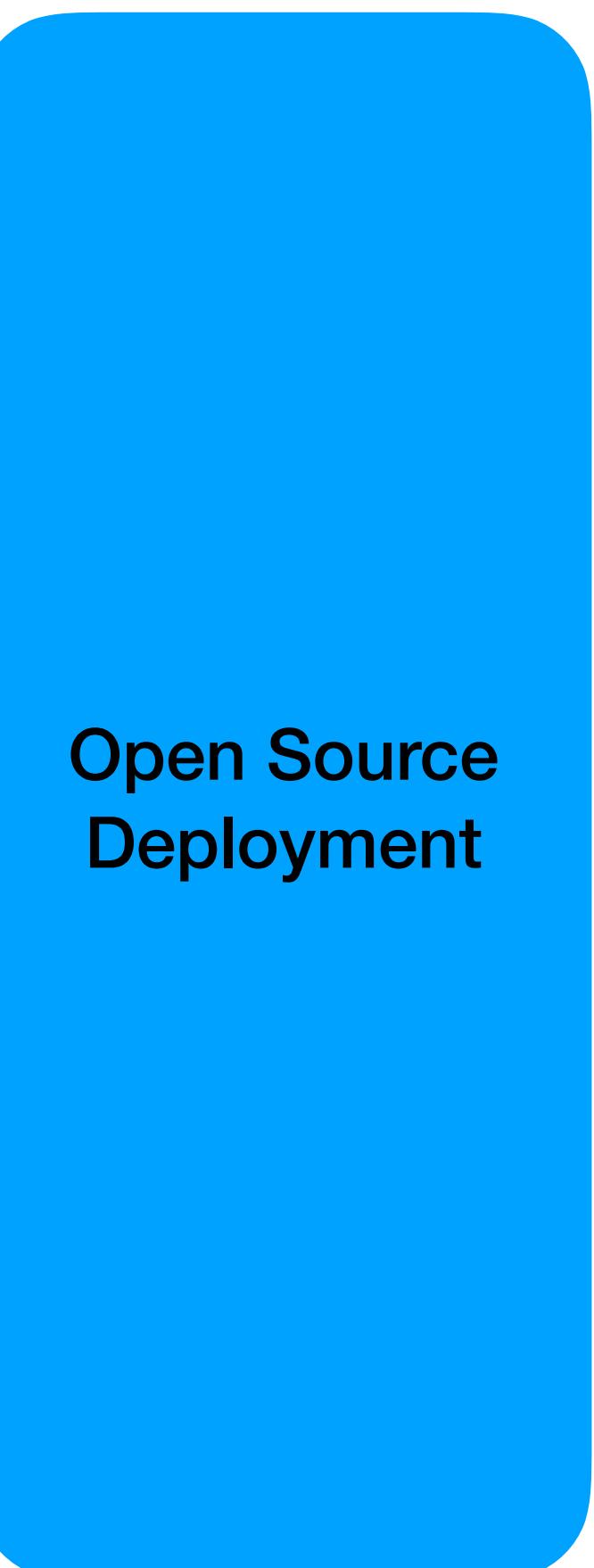
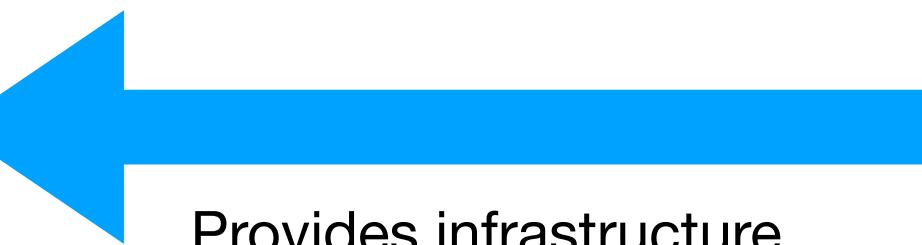
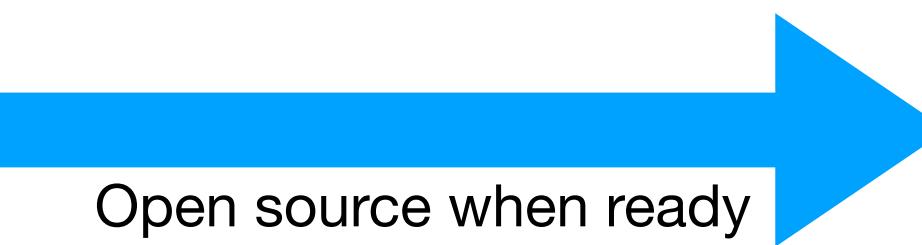
Compilers: Extensible support for future models, optimizations and hardware architectures

Systems: On-device and cloud-based training, distributed systems for ML

Computer Architecture: Extensible, energy efficient hardware designs for inference and training



- PL:** High-level support for future ML applications
- Compilers:** Extensible support for future models, optimizations and hardware architectures
- Systems:** On-device and cloud-based training, distributed systems for ML
- Computer Architecture:** Extensible, energy efficient hardware designs for inference and training



Open source compilers have transformed our industry



First major open source
compiler collection

Open source compilers have transformed our industry



First major open source
compiler collection

LLVM: Higher-Level IR, new
optimizations, easier extensibility



Open source compilers have transformed our industry

In the age of domain-specialized systems...



First major open source
compiler collection

LLVM: Higher-Level IR, new
optimizations, easier extensibility



Open source compilers have transformed our industry

In the age of domain-specialized systems...

Specialized compiler stack
for Deep Learning



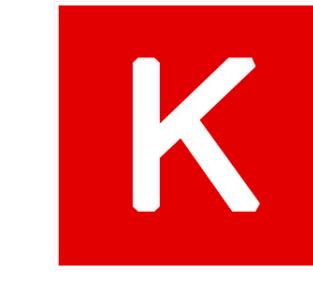
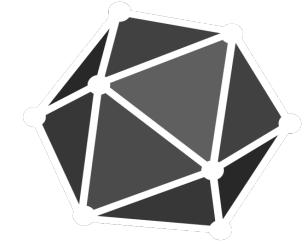
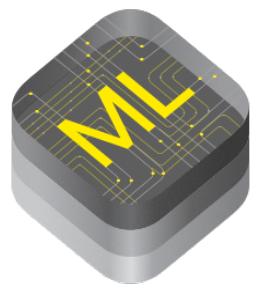
First major open source
compiler collection

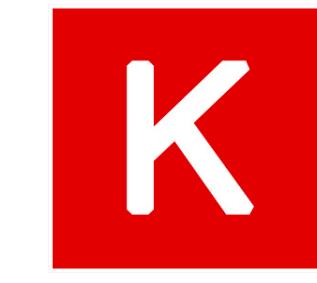
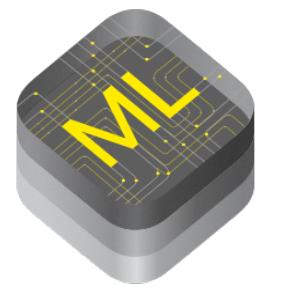
LLVM: Higher-Level IR, new
optimizations, easier extensibility



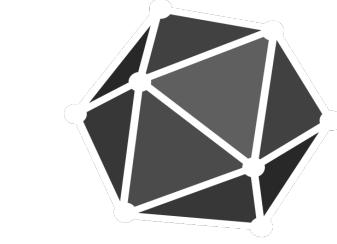
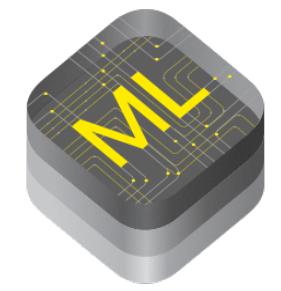
End the tyranny of closed deep learning systems!

Tianqi Chen



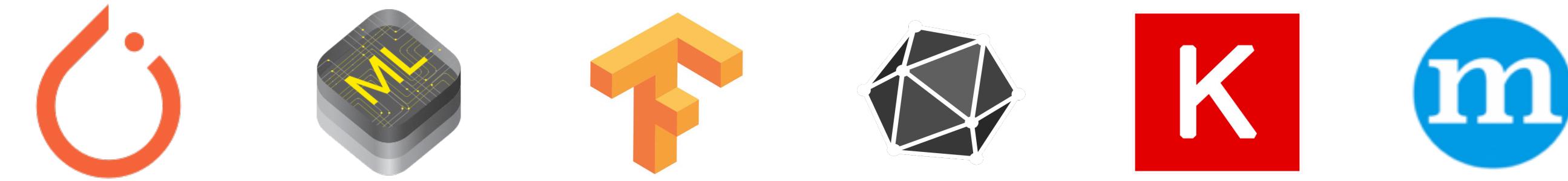


High-Level Differentiable IR



High-Level Differentiable IR

Tensor Expression IR

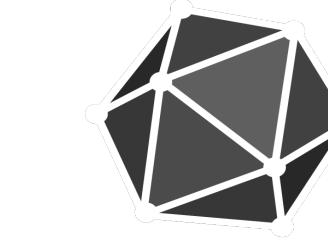
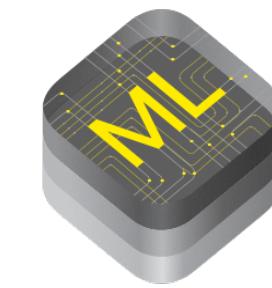


High-Level Differentiable IR

Tensor Expression IR

LLVM, CUDA, Metal





High-Level Differentiable IR

Tensor Expression IR

LLVM, CUDA, Metal

VTA



Edge
FPGA

Cloud
FPGA

ASIC



```
import tvm
from tvm import relay

graph, params =
    frontend.from_keras(keras_resnet50)
graph, lib, params =
    relay.build(graph, target)
```

Compile



```
import tvm
from tvm import relay

graph, params =
    frontend.from_keras(keras_resnet50)
graph, lib, params =
    relay.build(graph, target)
```

Compile



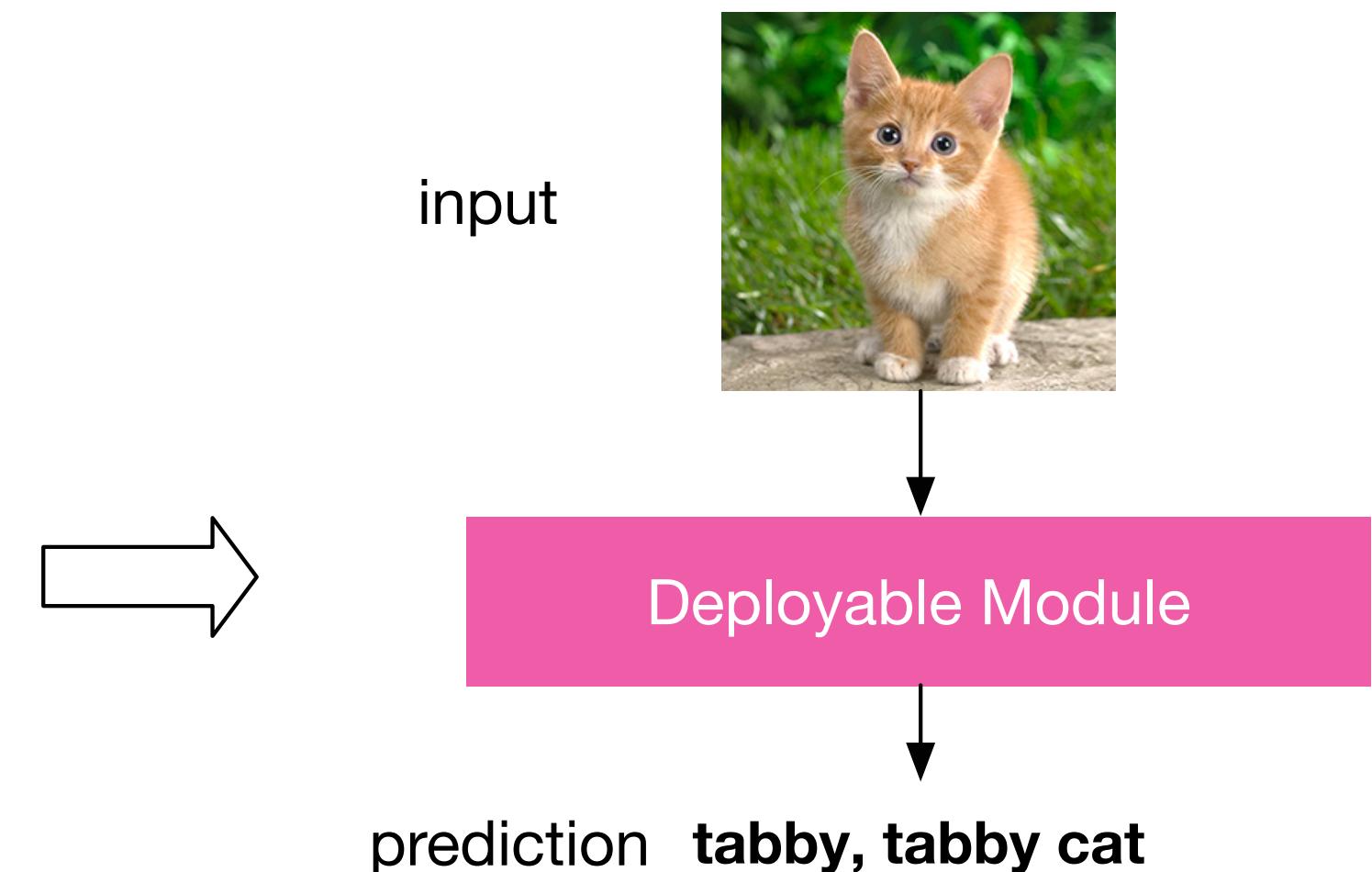


```
import tvm
from tvm import relay

graph, params =
    frontend.from_keras(keras_resnet50)
graph, lib, params =
    relay.build(graph, target)
```

Compile

```
module = runtime.create(graph, lib, tvm.gpu(0))
module.set_input(**params)
module.run(data=data_array)
output = tvm.nd.empty(out_shape, ctx=tvm.gpu(0))
module.get_output(0, output)
```



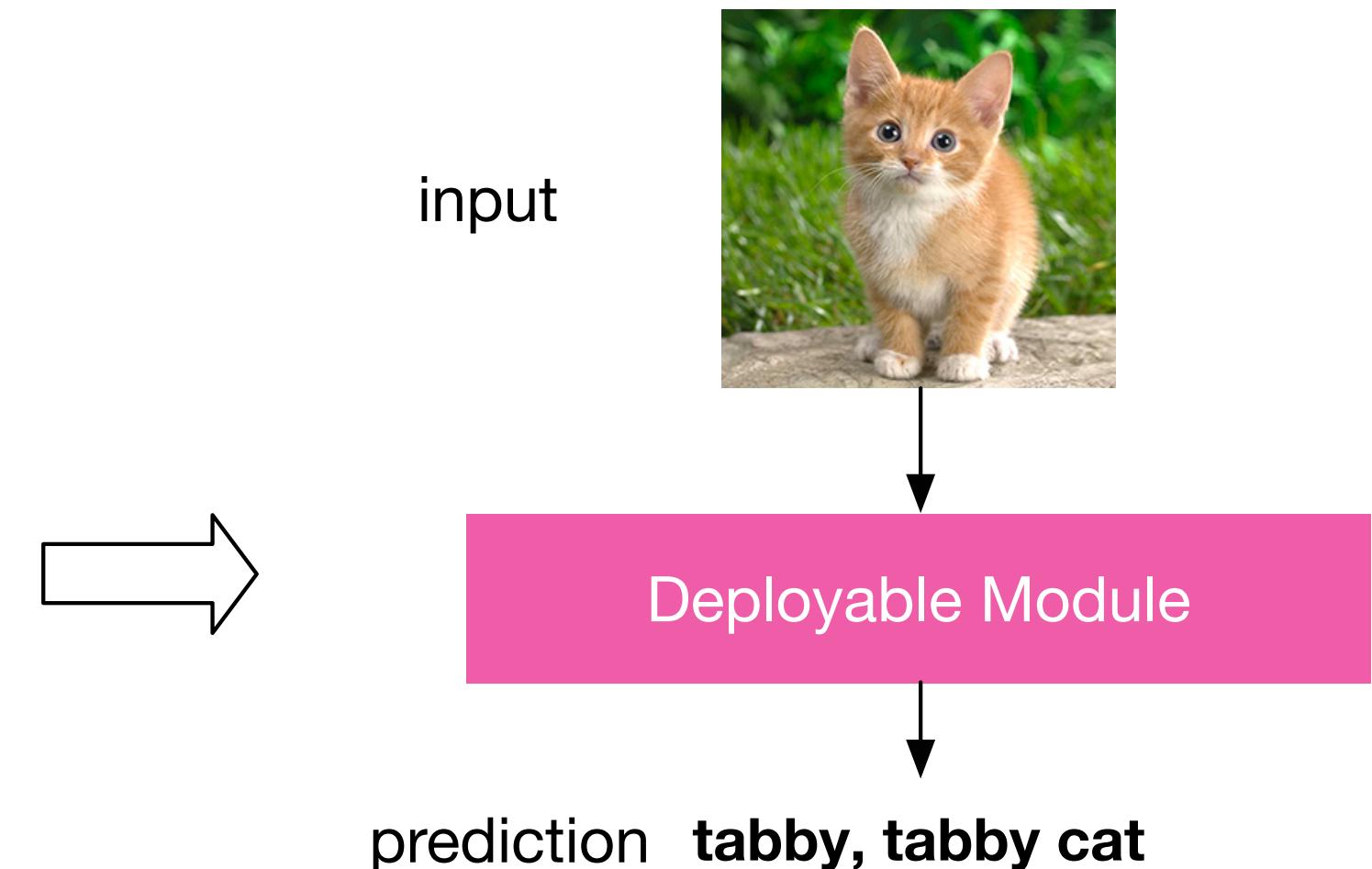


```
import tvm
from tvm import relay

graph, params =
    frontend.from_keras(keras_resnet50)
graph, lib, params =
    relay.build(graph, target)
```

Compile

```
module = runtime.create(graph, lib, tvm.gpu(0))
module.set_input(**params)
module.run(data=data_array)
output = tvm.nd.empty(out_shape, ctx=tvm.gpu(0))
module.get_output(0, output)
```

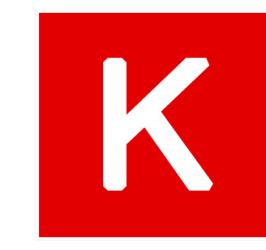
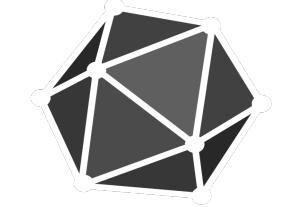
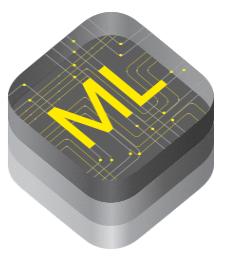


On languages and platforms you choose





Automated by Machine Learning



High-Level Differentiable IR

Tensor Expression IR

LLVM, CUDA, Metal

VTA



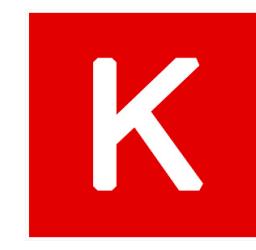
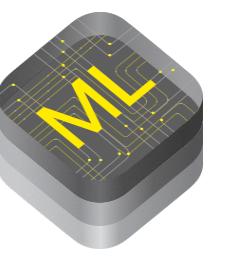
Edge
FPGA

Cloud
FPGA

ASIC



Automated by Machine Learning



High-Level Differentiable IR

Tensor Expression IR

LLVM, CUDA, Metal

VTA



Edge
FPGA

Cloud
FPGA

ASIC

Optimization

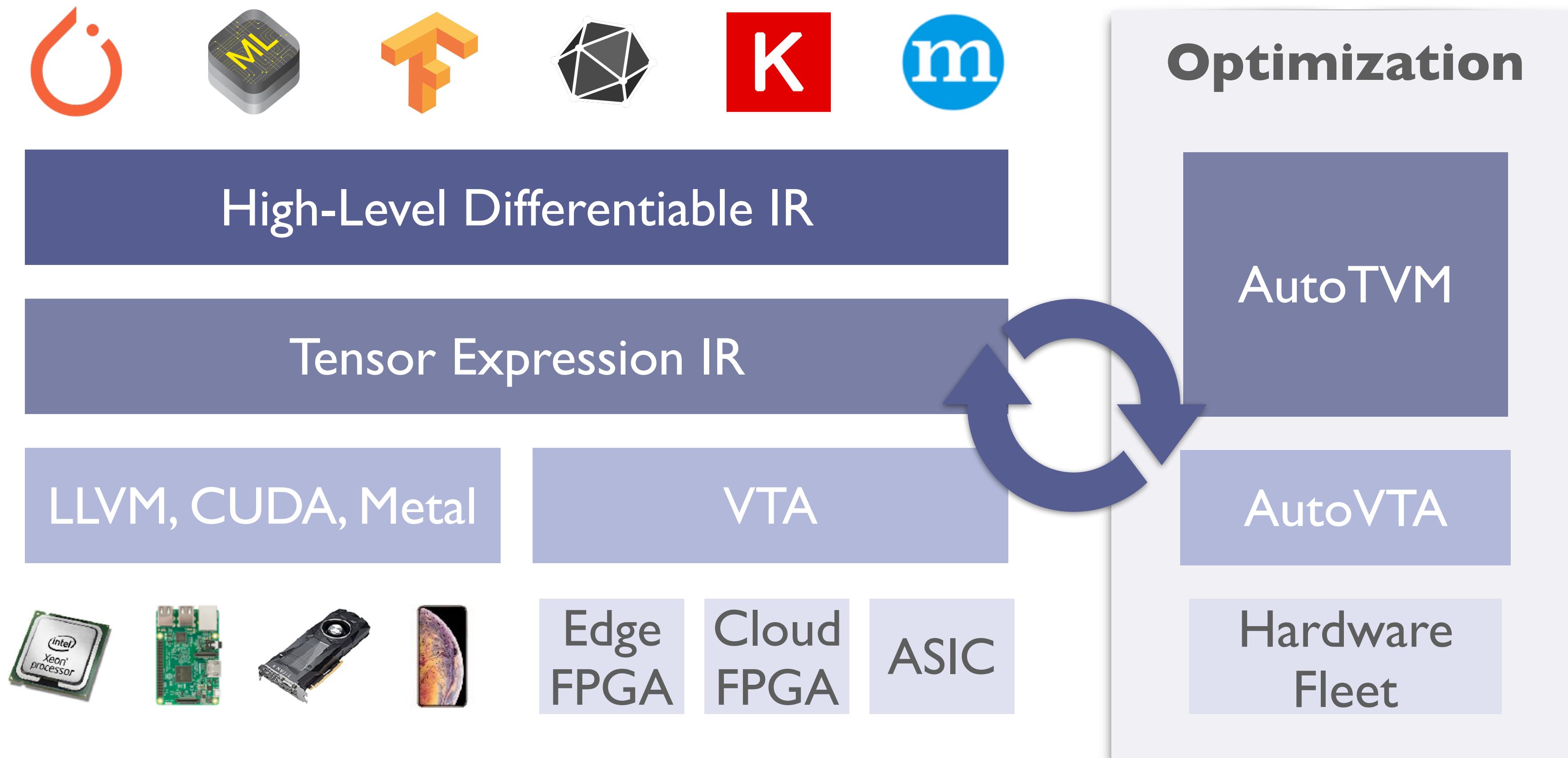
AutoTVM

AutoVTA

Hardware
Fleet

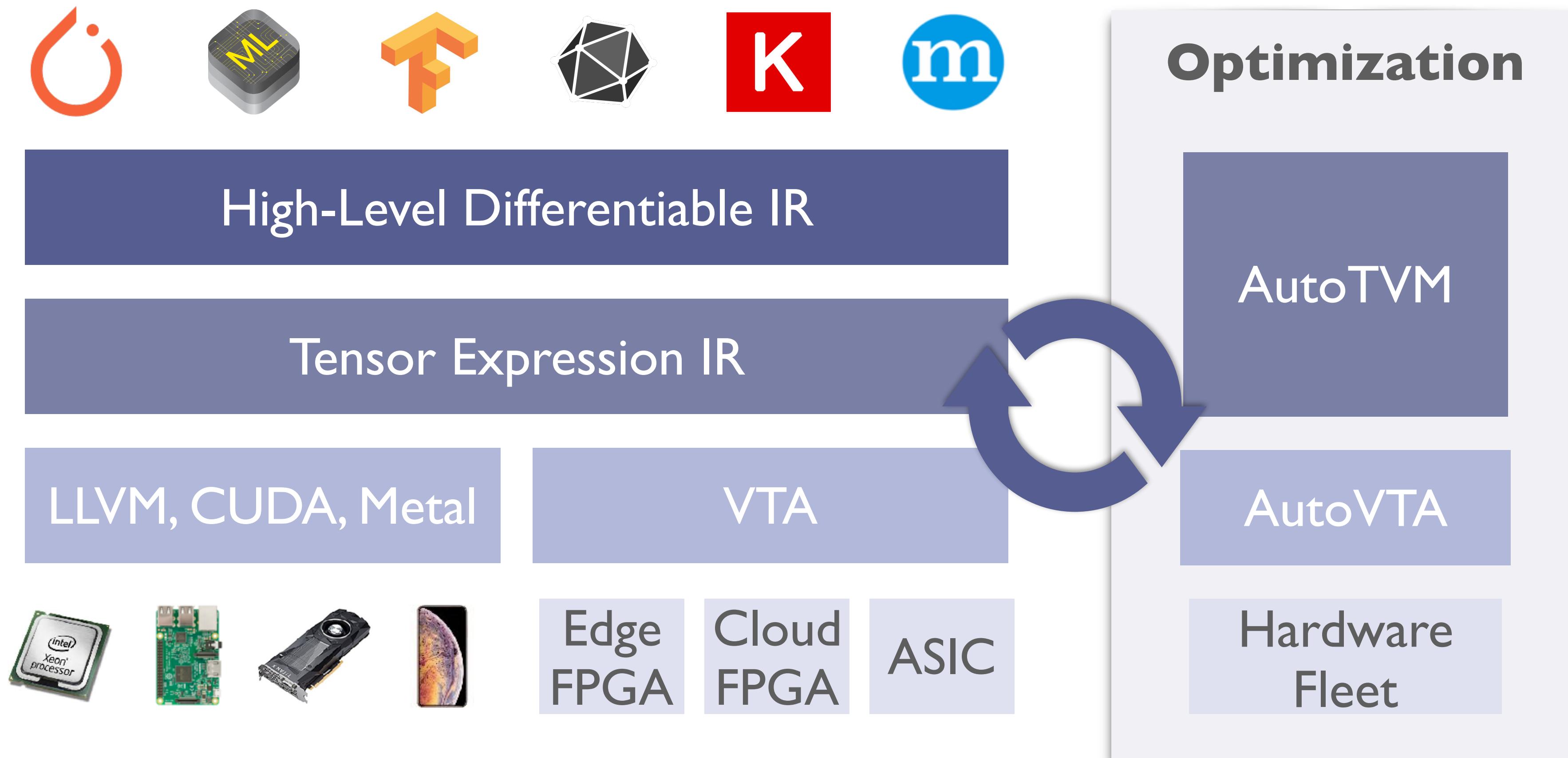


Automated by Machine Learning

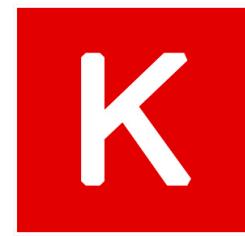
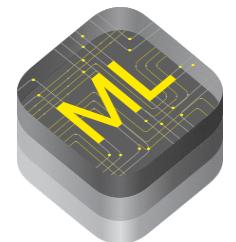




Automated by Machine Learning



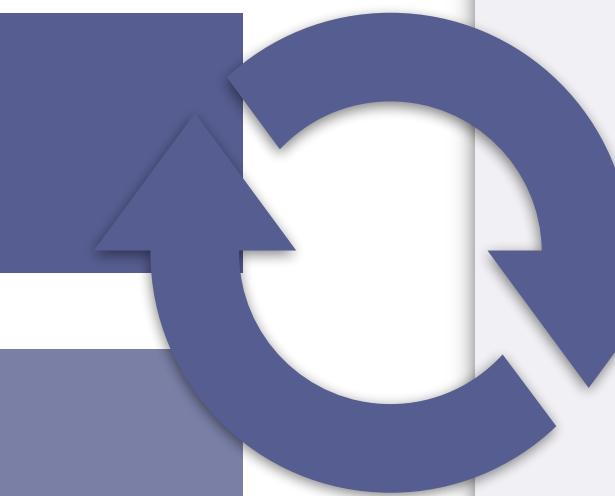
Diverse Hardware backends



Optimization

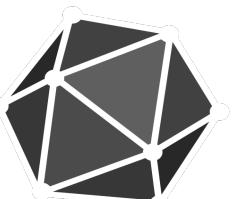
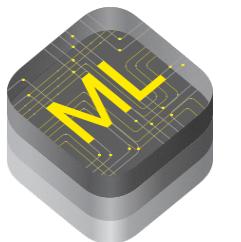
High-Level Differentiable IR

Tensor Expression IR



AutoTVM

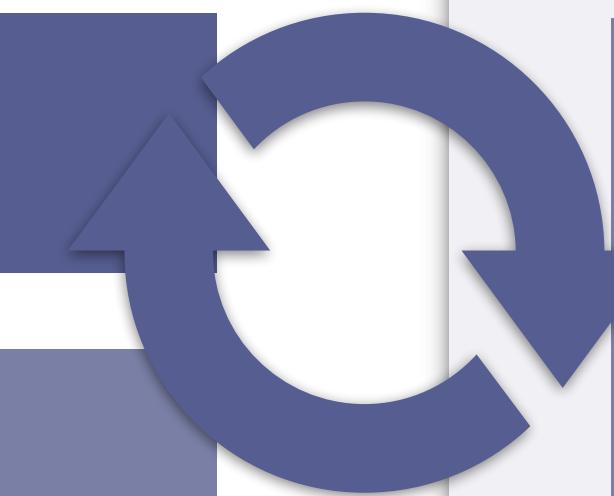
Diverse Hardware backends



Optimization

High-Level Differentiable IR

Tensor Expression IR



AutoTVM

LLVM

ARM

x86

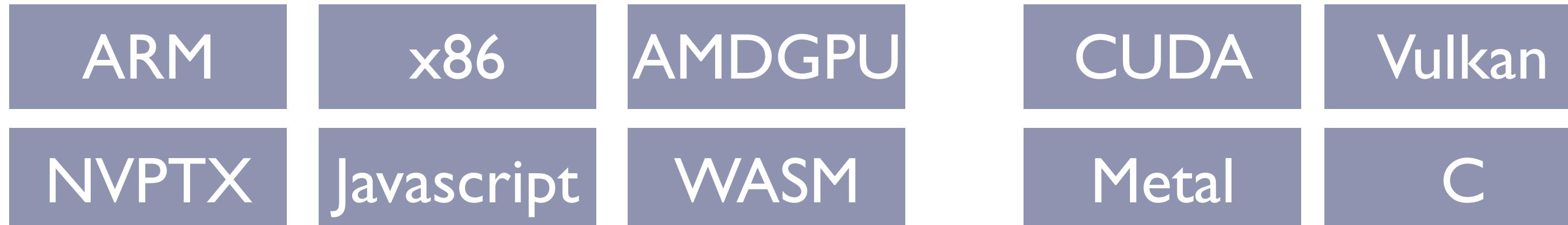
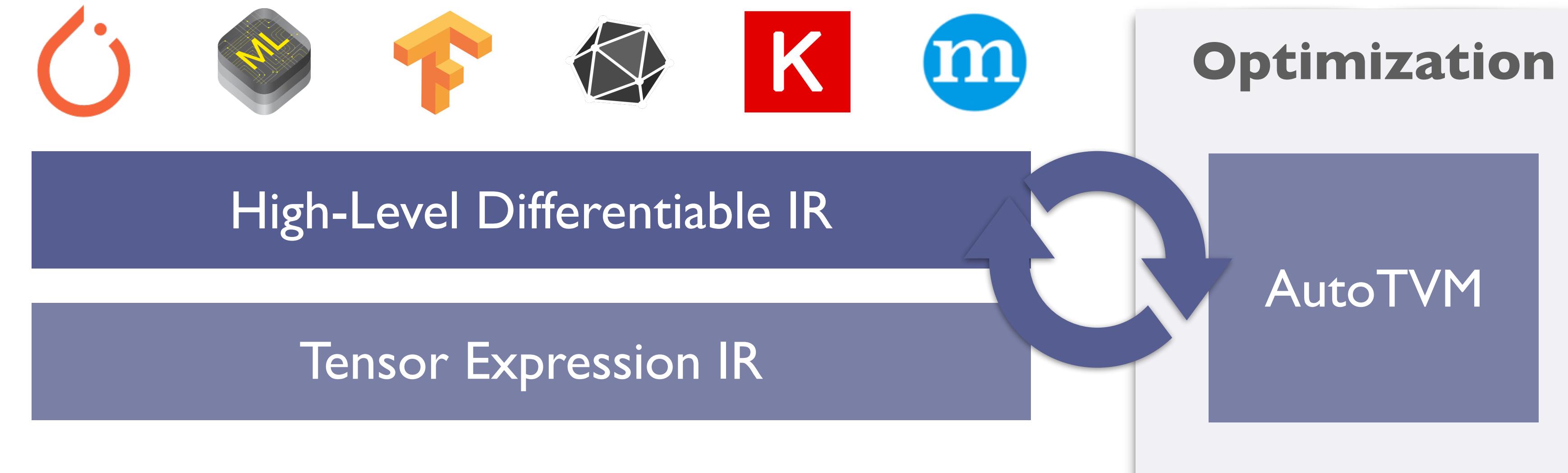
AMDGPU

NVPTX

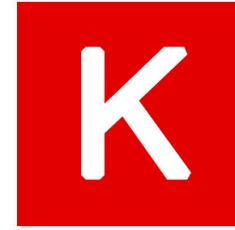
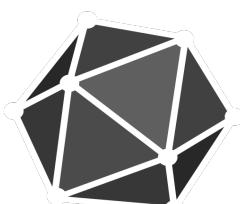
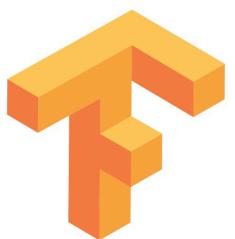
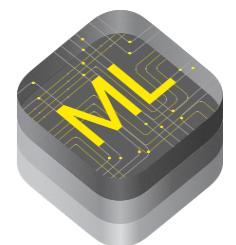
Javascript

WASM

Diverse Hardware backends



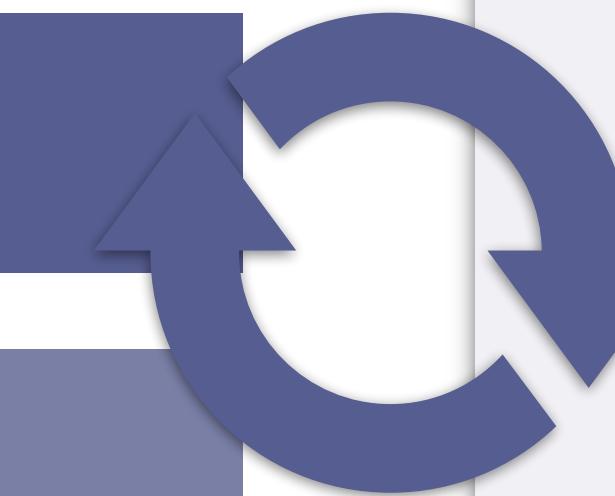
Diverse Hardware backends



Optimization

High-Level Differentiable IR

Tensor Expression IR



AutoTVM

LLVM

ARM

x86

AMDGPU

CUDA

Vulkan

VTA

NVPTX

Javascript

WASM

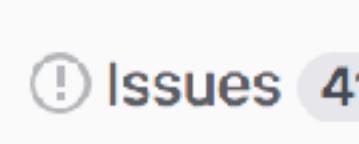
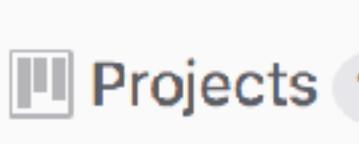
Metal

C

TVM Open Source Community

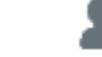
 [dmlc / tvm](#)

 [Unwatch](#) [283](#)  [Unstar](#) [2,449](#)  [Fork](#) [574](#)

 [Code](#)  [Issues 41](#)  [Pull requests 30](#)  [Projects 1](#)  [Wiki](#)  [Insights](#)  [Settings](#)

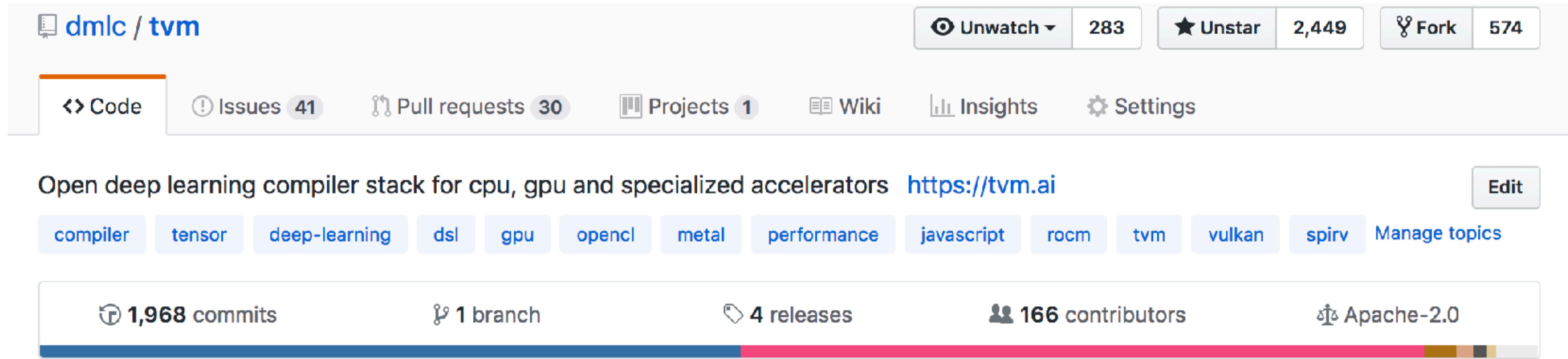
Open deep learning compiler stack for cpu, gpu and specialized accelerators <https://tvm.ai> 

[compiler](#) [tensor](#) [deep-learning](#) [dsl](#) [gpu](#) [opencl](#) [metal](#) [performance](#) [javascript](#) [rocm](#) [tvm](#) [vulkan](#) [spirv](#) [Manage topics](#)

 **1,968** commits  **1** branch  **4** releases  **166** contributors  **Apache-2.0**



TVM Open Source Community



The screenshot shows the GitHub repository page for `dmlc/tvm`. The page includes the following elements:

- Header:** `dmlc/tvm`, **Unwatch** (283), **Unstar** (2,449), **Fork** (574).
- Nav Bar:** **Code** (selected), **Issues** (41), **Pull requests** (30), **Projects** (1), **Wiki**, **Insights**, **Settings**.
- Description:** Open deep learning compiler stack for cpu, gpu and specialized accelerators <https://tvm.ai>.
- Topics:** compiler, tensor, deep-learning, dsl, gpu, opencl, metal, performance, javascript, rocm, tvm, vulkan, spirv, [Manage topics](#).
- Metrics:** 1,968 commits, 1 branch, 4 releases, 166 contributors, Apache-2.0.

Apache governance model: grant project ownership by merit.

11 committers, 29 reviewers, 166 contributors.

Contributed by the community, for the community.

Industrial Impact

TVM + AWS



Vin Sharma, Amazon SageMaker Neo

Amazon: vinarm@ | Twitter: [@ciphr](https://twitter.com/ciphr)

How is AWS using TVM?

How is AWS using TVM?

- As a back-end for Apache MXNet
 - To deploy easily onto edge devices
 - To improve performance on target hardware

How is AWS using TVM?

- As a back-end for Apache MXNet
 - To deploy easily onto edge devices
 - To improve performance on target hardware
- As an optimizer for Amazon AI services
 - Amazon Rekognition: To improve end-to-end latency
 - Amazon Alexa: To increase resource efficiency on Echo/Dot

How is AWS using TVM?

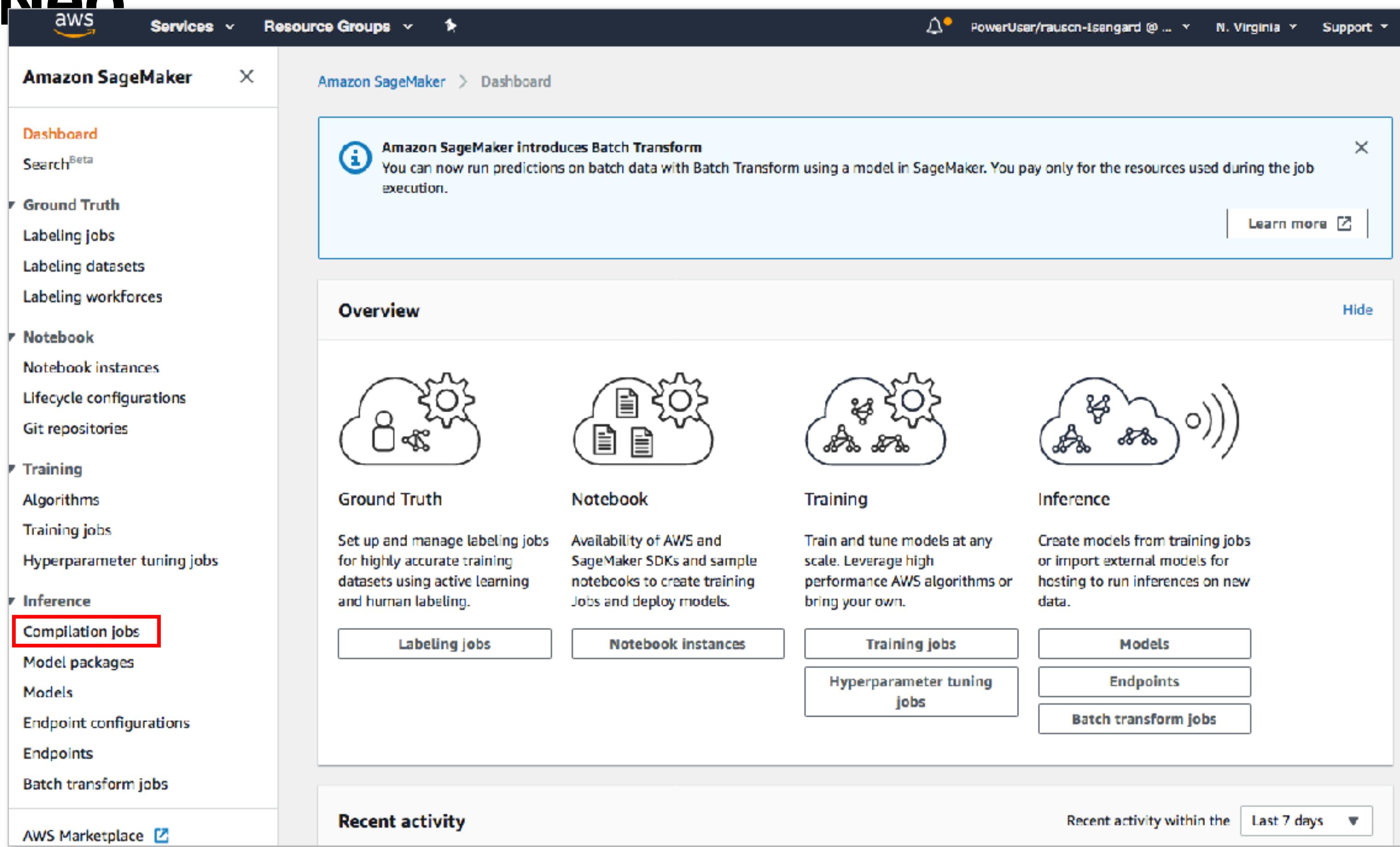
- As a back-end for Apache MXNet
 - To deploy easily onto edge devices
 - To improve performance on target hardware
- As an optimizer for Amazon AI services
 - Amazon Rekognition: To improve end-to-end latency
 - Amazon Alexa: To increase resource efficiency on Echo/Dot
- In a tool chain for Amazon Inferentia

How is AWS using TVM?

- As a back-end for Apache MXNet
 - To deploy easily onto edge devices
 - To improve performance on target hardware
- As an optimizer for Amazon AI services
 - Amazon Rekognition: To improve end-to-end latency
 - Amazon Alexa: To increase resource efficiency on Echo/Dot
- In a tool chain for Amazon Inferentia

How is AWS enabling adoption of TVM?

In a new service called **Amazon SageMaker Neo**



The screenshot shows the Amazon SageMaker Neo dashboard. The left sidebar has a red box around the 'Compilation jobs' link under the 'Inference' section. The main content area features a 'Batch Transform' announcement and an 'Overview' section with four categories: 'Ground Truth', 'Notebook', 'Training', and 'Inference', each with a description and associated icons.

Amazon SageMaker Neo Dashboard

Amazon SageMaker introduces Batch Transform

You can now run predictions on batch data with Batch Transform using a model in SageMaker. You pay only for the resources used during the job execution.

Overview

- Ground Truth**
Set up and manage labeling jobs for highly accurate training datasets using active learning and human labeling.
- Notebook**
Availability of AWS and SageMaker SDKs and sample notebooks to create training jobs and deploy models.
- Training**
Train and tune models at any scale. Leverage high performance AWS algorithms or bring your own.
- Inference**
Create models from training jobs or import external models for hosting to run inferences on new data.

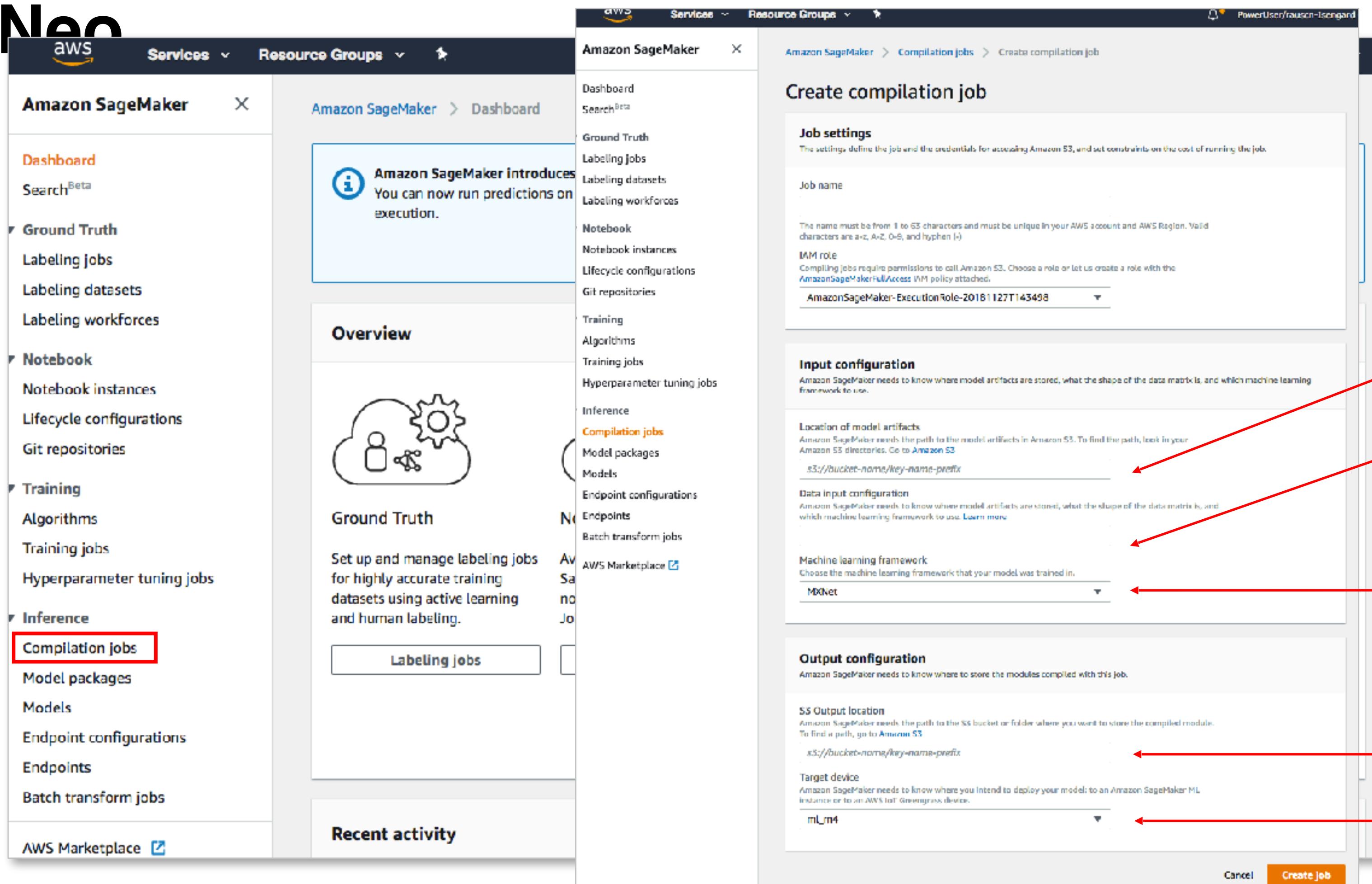
Recent activity

Recent activity within the **Last 7 days**

How is AWS enabling adoption of TVM?

In a new service called **Amazon SageMaker**

Neo



The screenshot shows the Amazon SageMaker console. The left sidebar navigation includes 'Dashboard', 'Ground Truth' (with 'Labeling jobs', 'Labeling datasets', 'Labeling workforces'), 'Notebook' (with 'Notebook instances', 'Lifecycle configurations', 'Git repositories'), 'Training' (with 'Algorithms', 'Training jobs', 'Hyperparameter tuning jobs'), and 'Inference' (with 'Compilation jobs' highlighted in red, 'Model packages', 'Models', 'Endpoint configurations', 'Endpoints', 'Batch transform jobs'). The main dashboard has a message: 'Amazon SageMaker introduces You can now run predictions on execution.' Below it is an 'Overview' section with a cloud icon and text about setting up labeling jobs for training datasets. The right panel is titled 'Create compilation job' and contains several configuration sections: 'Job settings' (Job name: 'MyCompilationJob', IAM role: 'AmazonSageMaker-ExecutionRole-20181127T143498'), 'Input configuration' (Location of model artifacts: 's3://bucket-name/key-name-prefix'), 'Data input configuration' (Location of data artifacts: 's3://bucket-name/key-name-prefix'), 'Machine learning framework' (MXNet), 'Output configuration' (SS Output location: 's3://bucket-name/key-name-prefix'), and 'Target device' (mLm4). Red arrows point from the text labels on the right to the corresponding fields in the 'Input configuration', 'Machine learning framework', 'Output location', and 'Target Platform' sections of the dialog.

Model input files:
MXNet: .json & .params

Name and shape of input node:
{“data”:[1,3,227,277]}

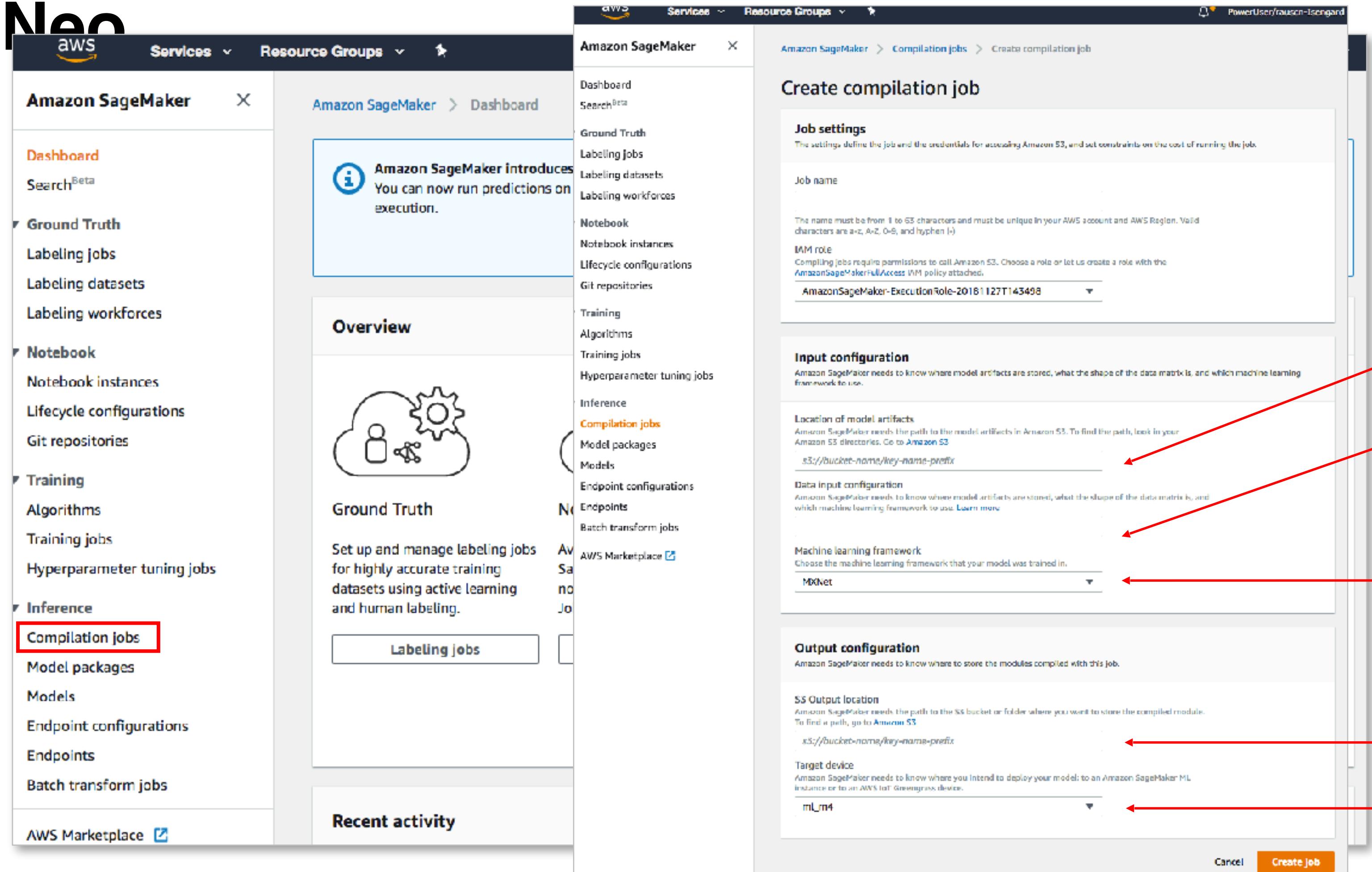
Framework

Output Location

Target Platform:
Cloud Instance Type I Edge Device

How is AWS enabling adoption of TVM?

In a new service called **Amazon SageMaker Neo**



The screenshot shows the Amazon SageMaker Neo interface. On the left is a navigation sidebar with sections like Dashboard, Ground Truth, Labeling jobs, Labeling datasets, Labeling workforces, Notebook, Training, Inference, and a highlighted 'Compilation jobs' section. The main dashboard shows an 'Amazon SageMaker introduces' message and an 'Overview' section with a cloud icon. The right side shows the 'Create compilation job' dialog, which includes fields for 'Job settings' (Job name: 'AmazonSageMaker-ExecutionRole-20181127T143498'), 'Input configuration' (Location of model artifacts: 's3://bucket-name/key-name-prefix'), 'Machine learning framework' (MXNet), 'Output configuration' (S3 Output location: 's3://bucket-name/key-name-prefix'), and 'Target device' (ml.m4). Red arrows point from the text labels on the right to the corresponding fields in the dialog.

Model input files:
MXNet: .json & .params

Name and shape of input node:
{“data”:[1,3,227,277]}

Framework

Output Location

Target Platform:
Cloud Instance Type | Edge Device

How is AWS contributing to TVM?

Releasing all TVM modifications and enhancements in Neo to open source

- **Frameworks:** TensorFlow, MXNet, PyTorch, ONNX
- **Models:** ResNet, VGG, Inception, MobileNet, DenseNet, SqueezeNet
- **Operators:** Several new ops in NNVM/TVM
- **Optimizations:** Node Annotation, Graph Partitioning, Ring Buffer, NHWC, Graph Tuning
- **Acceleration Library:** Nvidia TensorRT
- **Hardware:** Cross-Compilation to ARM, Intel, Nvidia; More Coming Soon

How is AWS contributing to TVM?

Releasing all TVM modifications and enhancements in Neo to open source

- **Frameworks:** TensorFlow, MXNet, PyTorch, ONNX
- **Models:** ResNet, VGG, Inception, MobileNet, DenseNet, SqueezeNet
- **Operators:** Several new ops in NNVM/TVM
- **Optimizations:** Node Annotation, Graph Partitioning, Ring Buffer, NHWC, Graph Tuning
- **Acceleration Library:** Nvidia TensorRT
- **Hardware:** Cross-Compilation to ARM, Intel, Nvidia; More Coming Soon

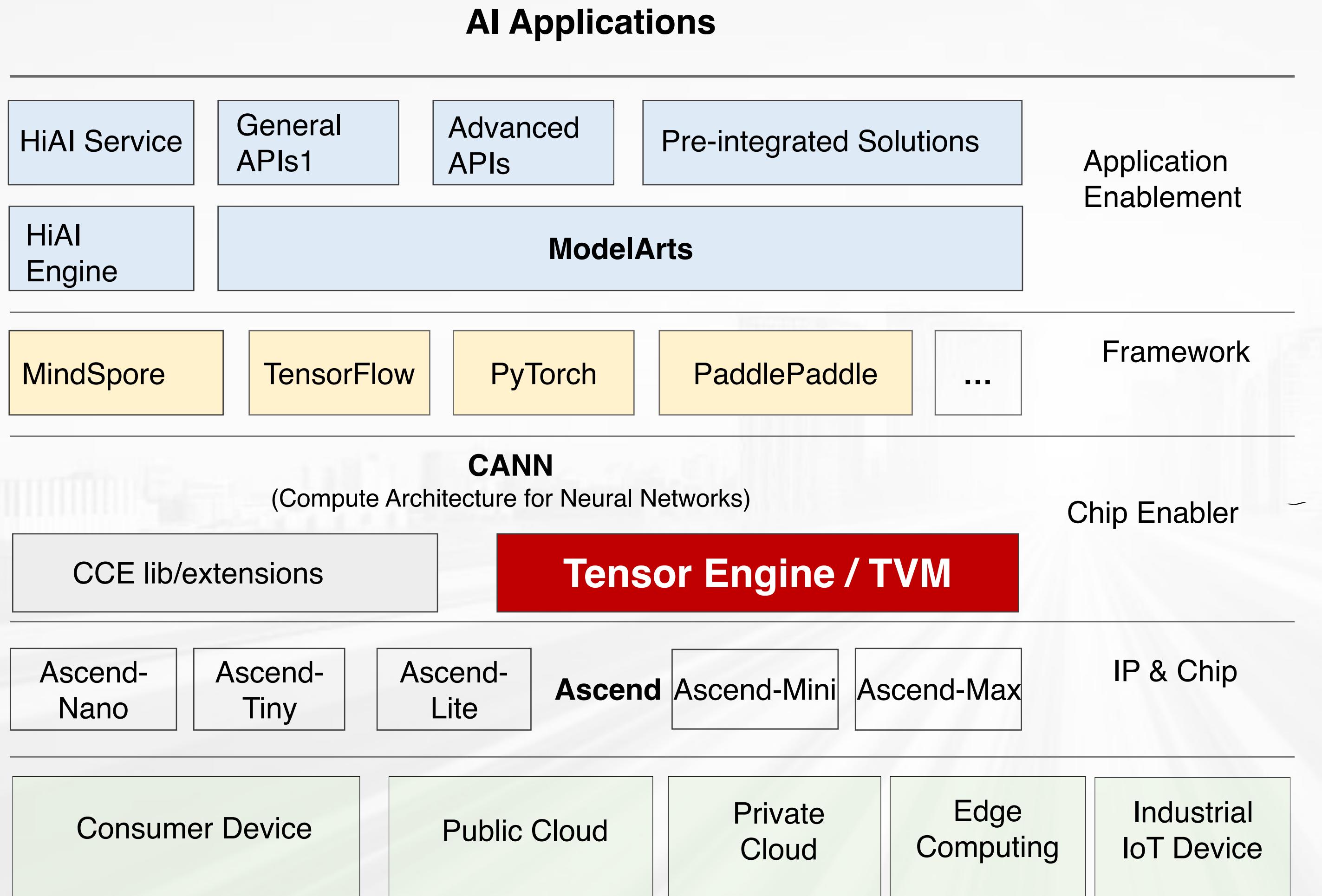


W PAUL G. ALLEN SCHOOL
OF COMPUTER SCIENCE & ENGINEERING



Chen Tian, Technical VP

TVM on Huawei's AI portfolio



Application enabling:

Full-pipeline services(ModelArts), hierarchical APIs, and pre-integrated solutions

MindSpore:

Unified training and inference framework for device, edge, and cloud (both standalone and cooperative)

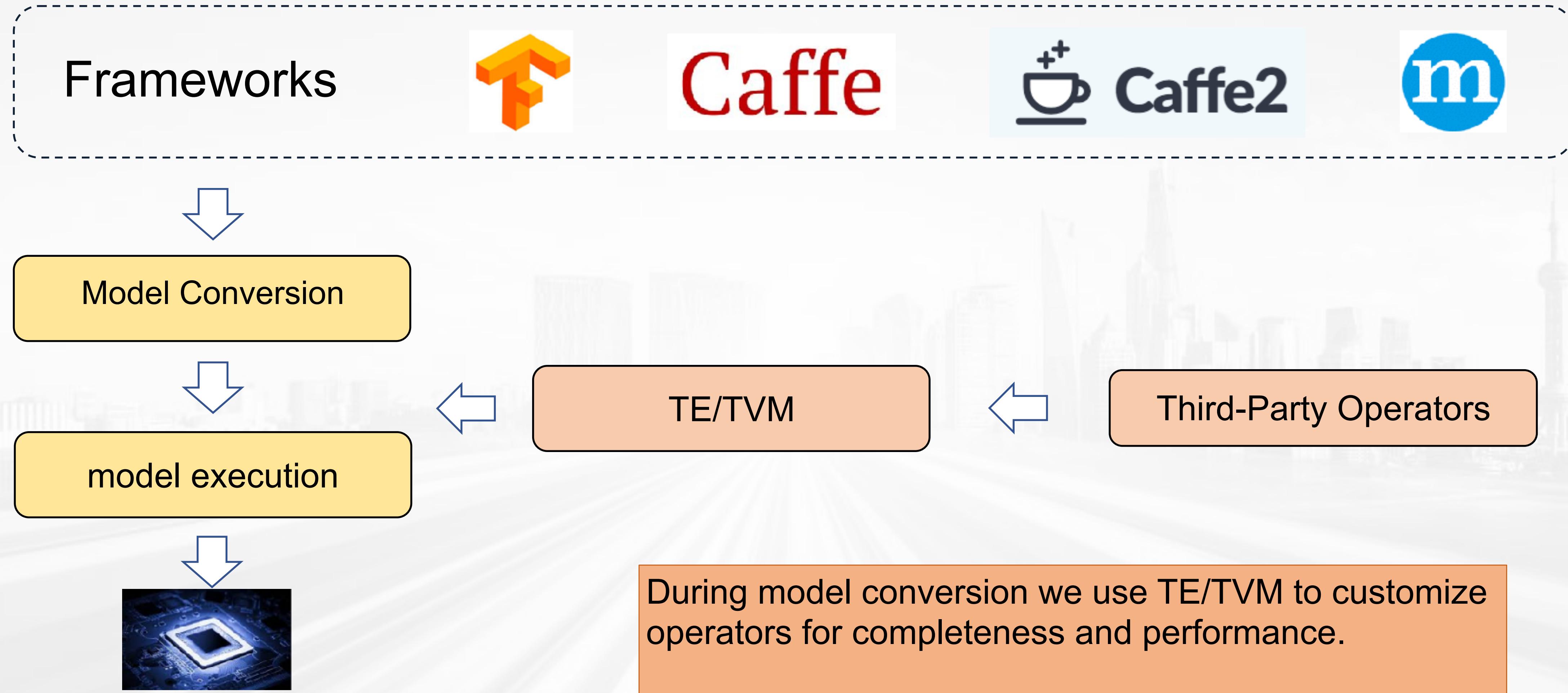
CANN:

Chip operators library and highly automated operators development toolkit

Ascend:

AI chip series based on unified scalable architecture

How do we use TVM



70+ operators are written by TVM , bring us ~3x development efficiency improvement

Successful Practice with Audi in Level 4 Autonomous Driving

~ A Complete City Commute Record ~

Driving in the evening



High-speed cruise



Traffic Jam Pilot (TJP)



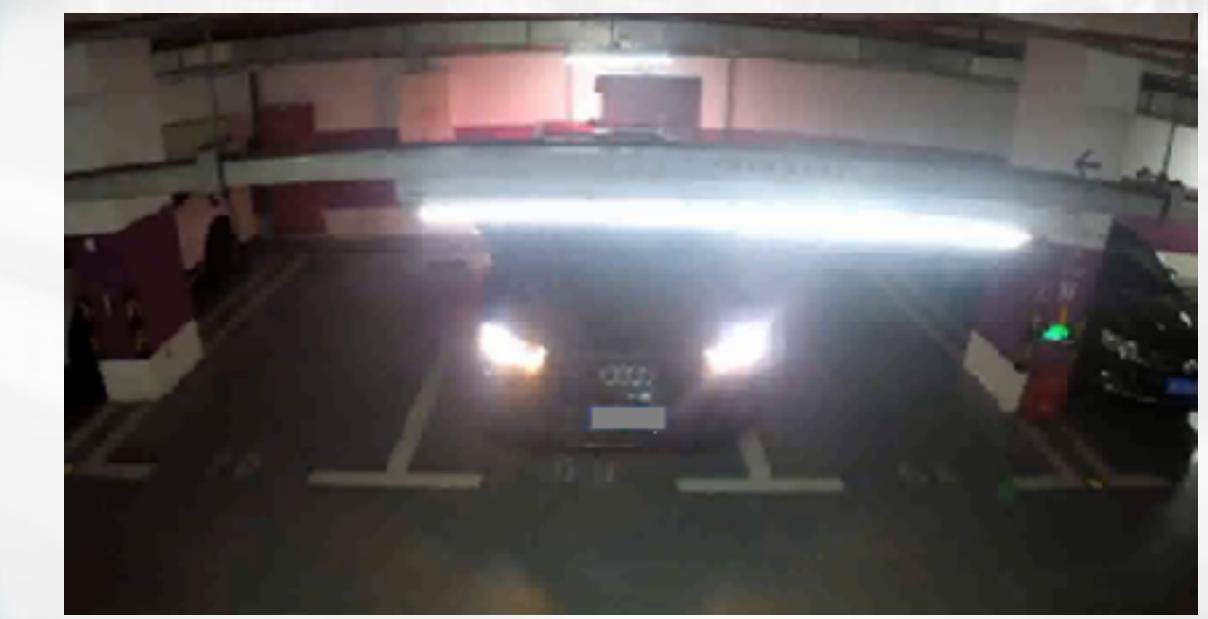
Traffic light identification



Pedestrian identification



Automatic parking



Joint developed autonomous driving algorithm gains leading scores in industry authoritative KITTI 2D/3D/BEV tests!

TVM is working on Atlas series product

Atlas 200 Developer Kit



- 16 TOPS INT8@24 W
- 1 USB type-C, 2 CCM interfaces, 1 GE network port, 1 SD card slot
- 8 GB memory

Atlas 300 AI Accelerator Card



- 64 TOPS INT8@75 W
- 64-channel HD video real-time analysis and JPEG decoding
- 32 GB memory, 204.8 GB/s memory bandwidth
- PCIe 3.0 x16, half-height half-length card

Atlas 500 AI Edge Station



- Capable of processing 16-channel HD videos in the size of a set-top-box (STB)
- Delivers 4x higher performance over counterparts

Atlas 800 AI Appliance



- Provides optimized AI environment based on the standard framework and programming environment
- Leverages high-performance GPU scheduling algorithms, improving resource utilization by over 15%

Smart Manufacturing

(intelligent quality inspection and flexible manufacturing)



Intelligent Care

(kindergarten and elderly care)



Smart Transportation

(traffic light tuning, intelligent traffic guiding)



Huawei's Contributions on TVM

8 Contributors:

kun-zh, sgrechanik-h, libing4752, derisavi-huawei, solin319, ehsanmok, gaoxiong-1, jiacunjiang1215

4 Reviewers:

Srkreddy1238 , PariksheetPinjari909 , siju-Samuel , Xqdan

We are working on:

1. Huawei Ascend ASIC support.
2. Front end to support Darknet, ONNX.
3. Optimization on Auto-TVM, IR extensions.
4. Tensorize, cache read/write, access_ptr API.

In the future we will try to:

1. Codegen for fused operators.
2. NLP support.
3. More optimization.
4. Training Operators.

Meghan Cowan

VGG11 on Raspberry Pi 3B



TensorflowLite

32bit fp

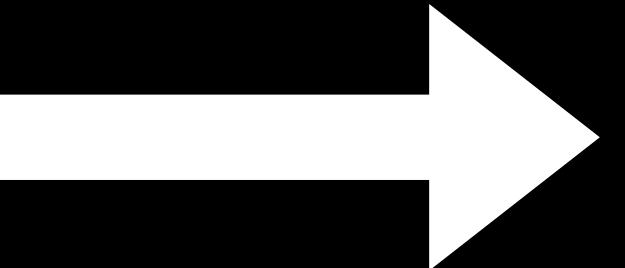
66% top-1 ImageNet accuracy

1.42 fps

VGG11 on Raspberry Pi 3B



Trained
binarized model



Operators
implemented with TVM

TensorflowLite
32bit fp

66% top-1 ImageNet accuracy
1.42 fps

VGG11 on Raspberry Pi 3B



Trained
binarized model



Operators
implemented with TVM

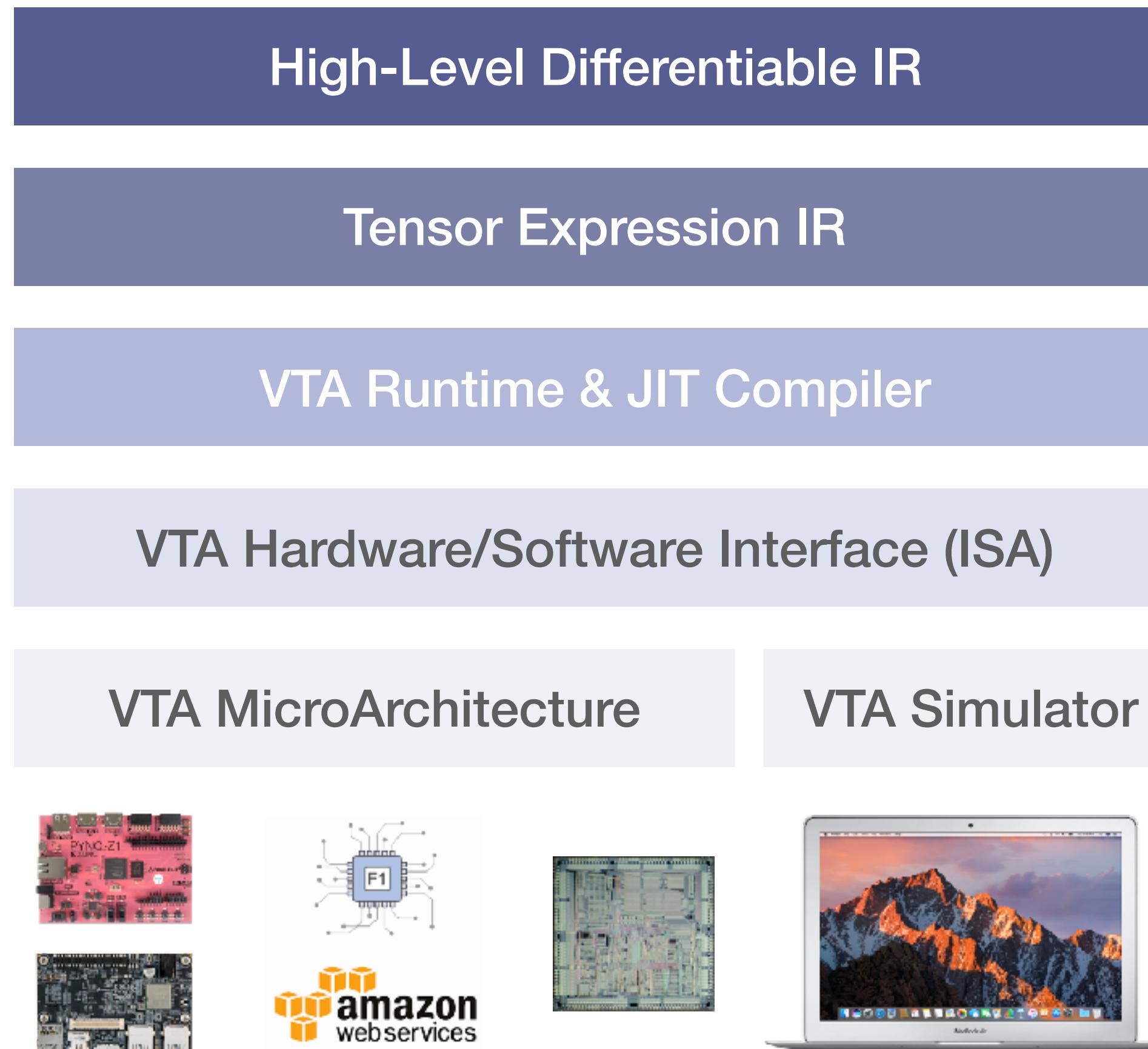
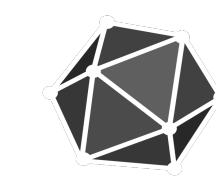
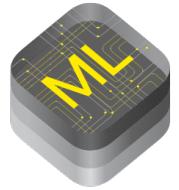
TensorflowLite
32bit fp
66% top-1 ImageNet accuracy
1.42 fps

TVM
2-bit activation 1-bit weight
62% top-1 ImageNet accuracy
4.67 fps

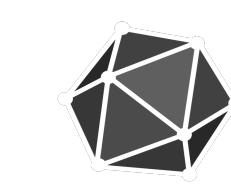
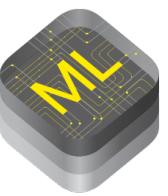
Further down the stack...

Thierry Moreau

Open Source Stack Overview



Open Source Stack Overview



VTA Backends

- **Simulator:** out-of-the-box testing to write compiler passes

Versatile Tenso Accelerator Stack (VTA)

High-Level Differentiable IR

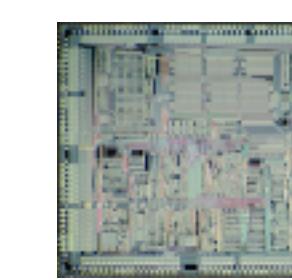
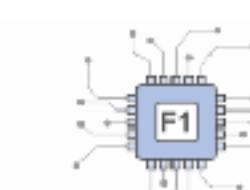
Tensor Expression IR

VTA Runtime & JIT Compiler

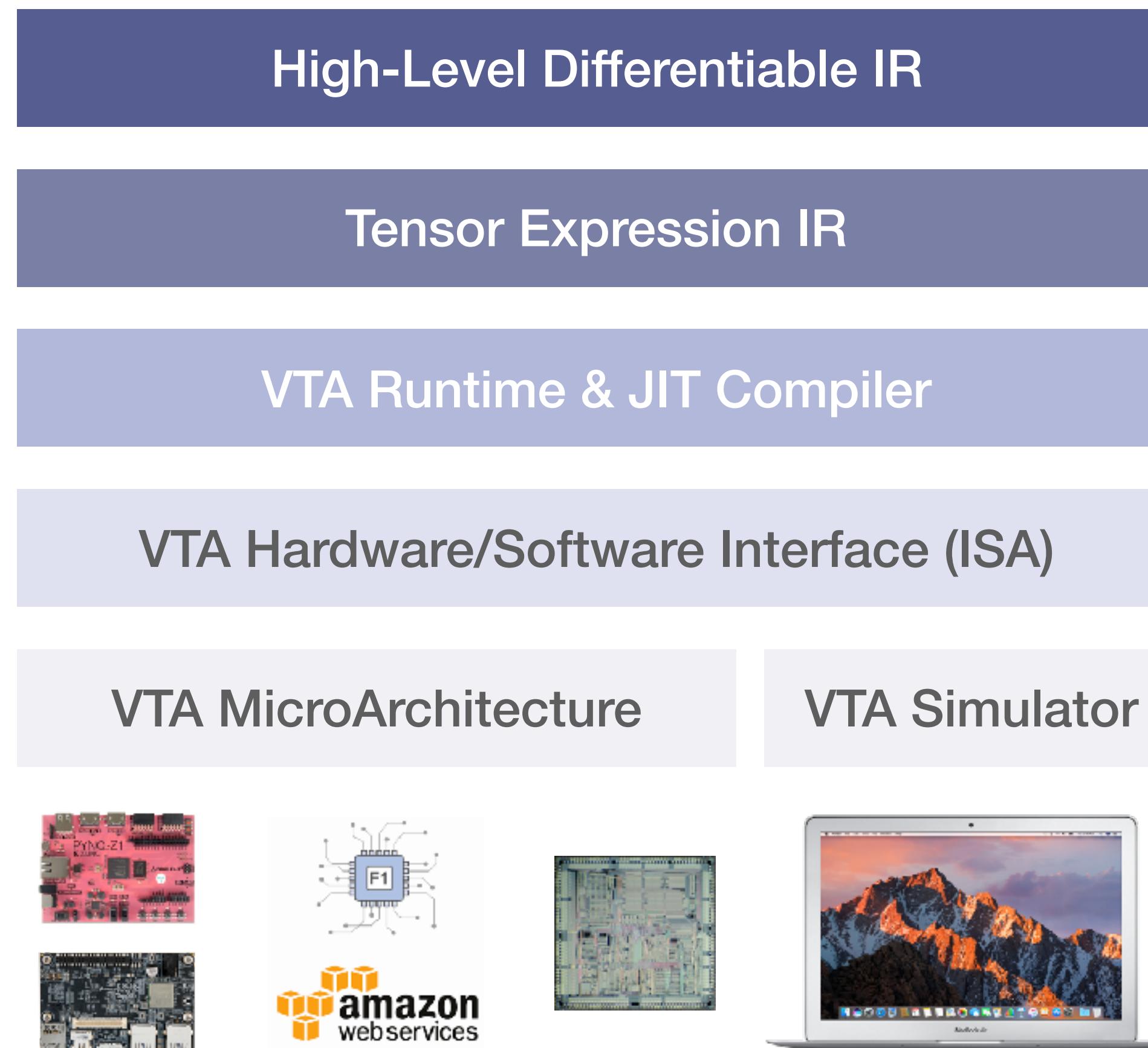
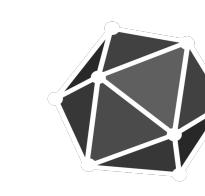
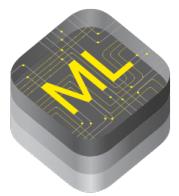
VTA Hardware/Software Interface (ISA)

VTA MicroArchitecture

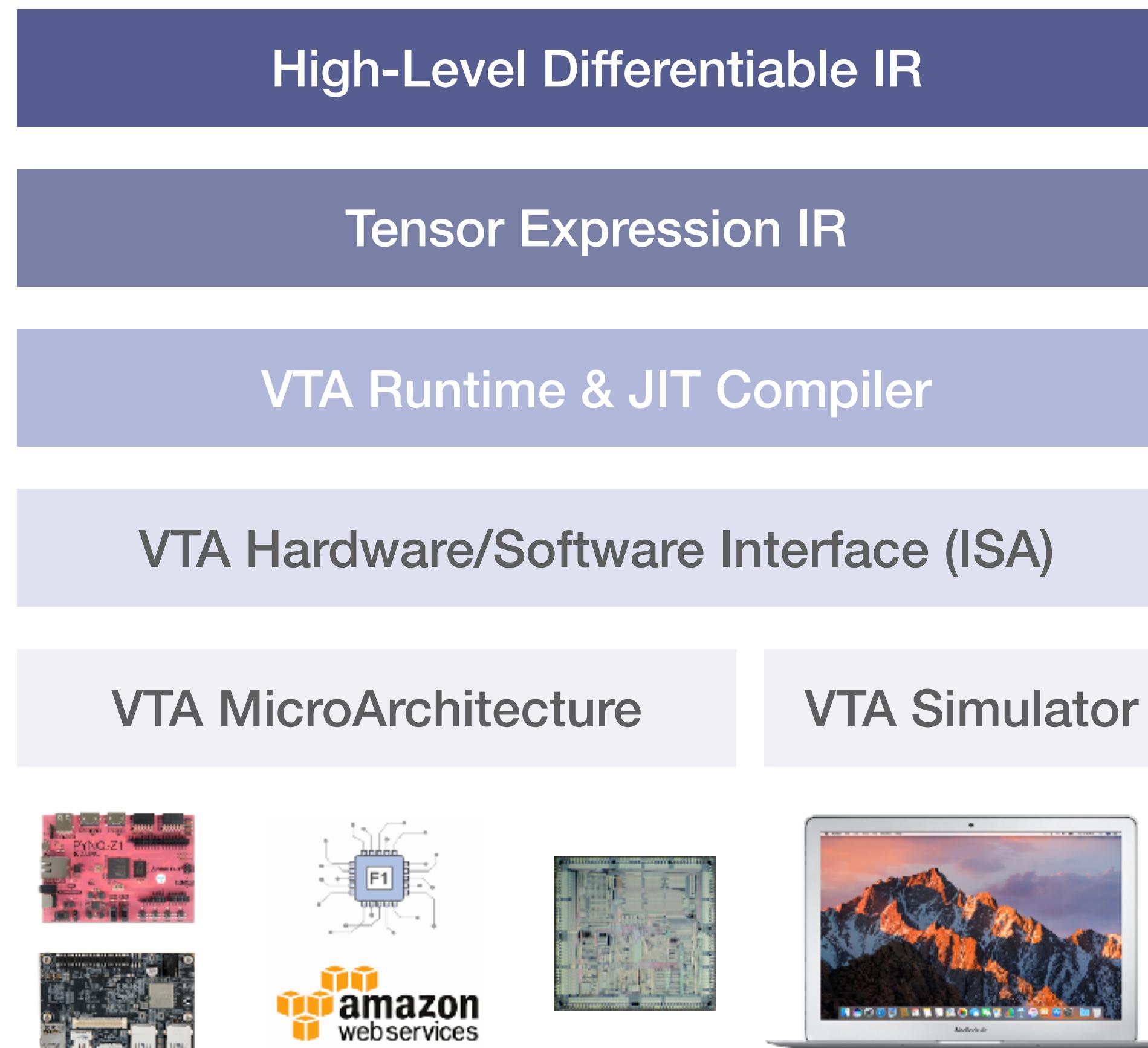
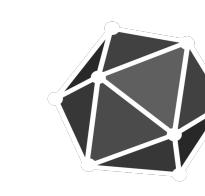
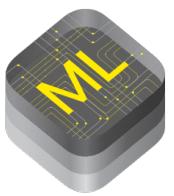
VTA Simulator



Open Source Stack Overview



Open Source Stack Overview



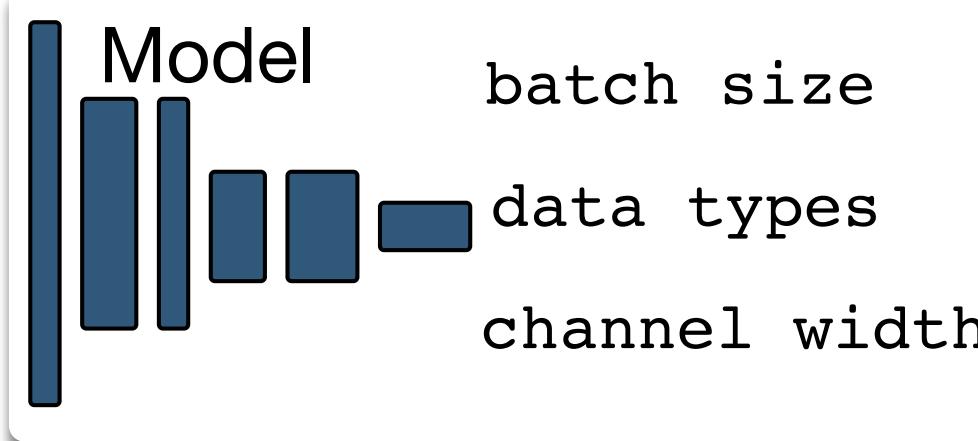
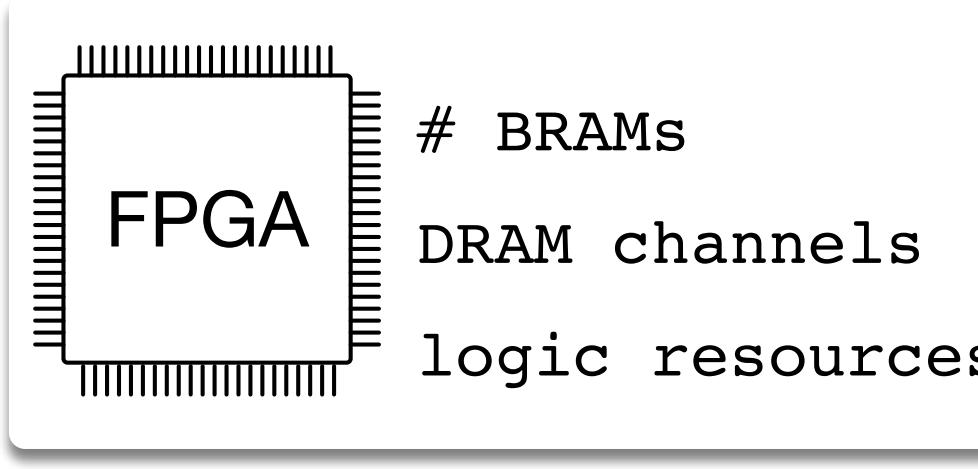
VTA Backends

- **Simulator:** out-of-the-box testing to write compiler passes
- **FPGA:** fast design iteration, quick deployment, flexibility
- **ASIC:** industrial-strength efficiency

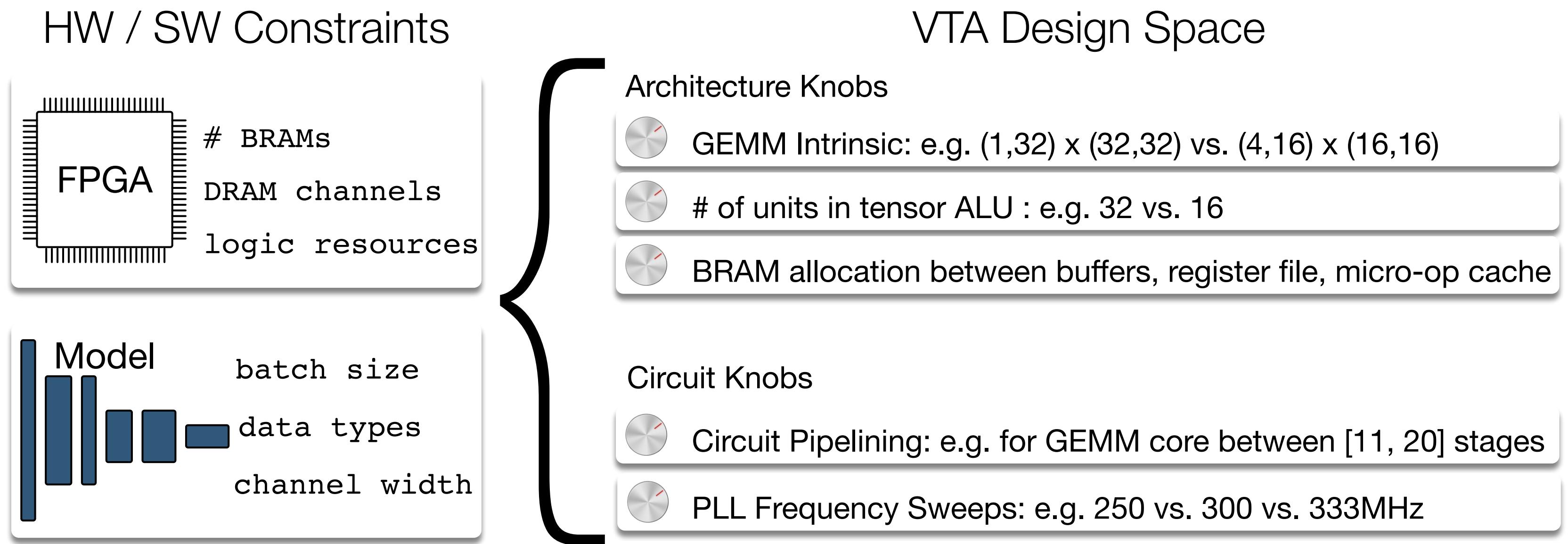
Versatile Tensor Accelerator Stack (VTA)

Hardware Exploration with VTA

HW / SW Constraints

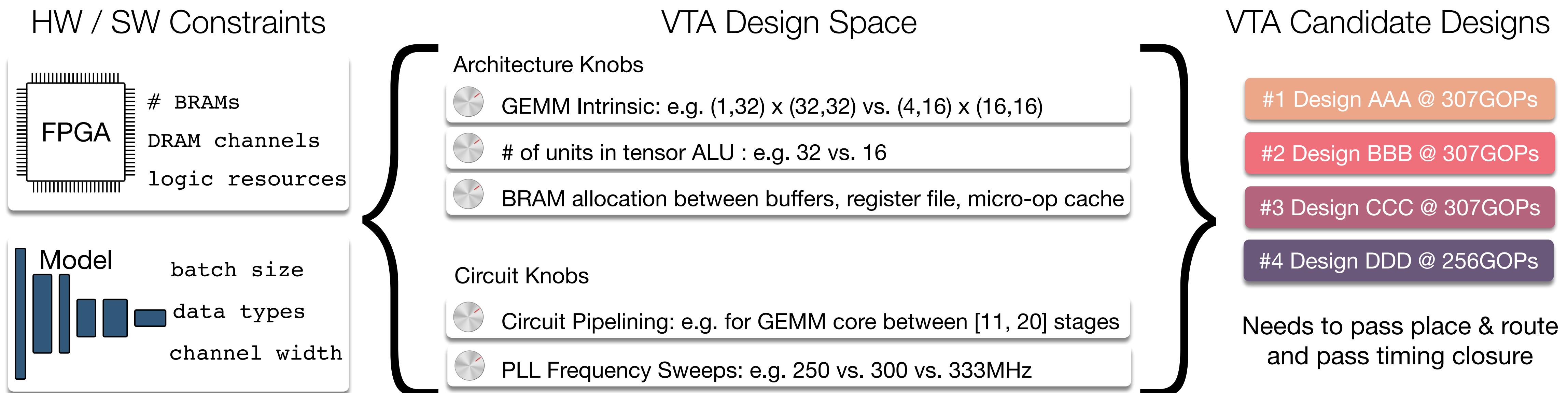


Hardware Exploration with VTA



1000s

Hardware Exploration with VTA



1000s

~ 10

Schedule Exploration with VTA

VTA Candidate Designs

#1 Design AAA @ 307GOPs

#2 Design BBB @ 307GOPs

#3 Design CCC @ 307GOPs

#4 Design DDD @ 256GOPs

Needs to pass place & route
and pass timing closure

Schedule Exploration with VTA

VTA Candidate Designs

#1 Design AAA @ 307GOPs

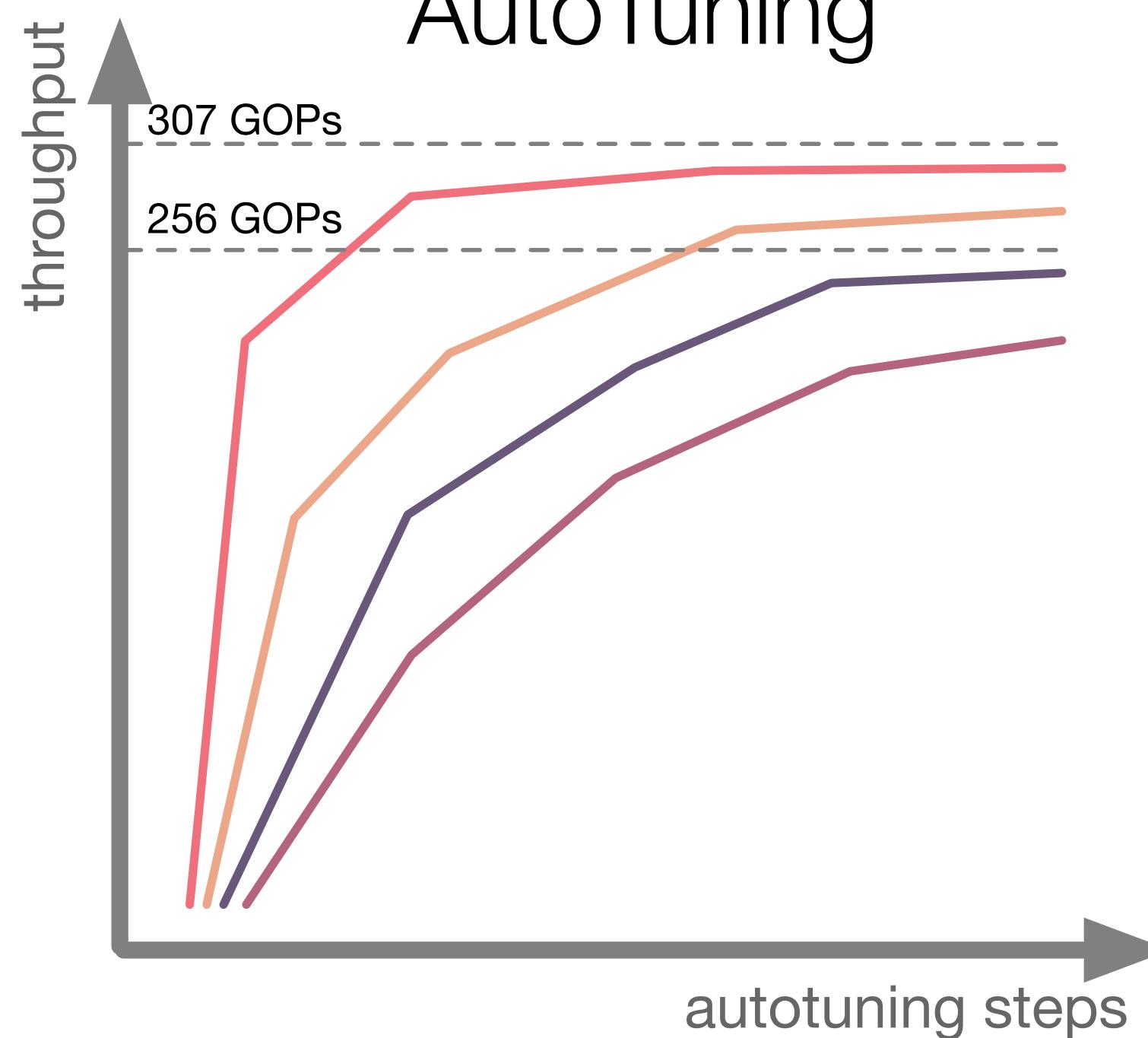
#2 Design BBB @ 307GOPs

#3 Design CCC @ 307GOPs

#4 Design DDD @ 256GOPs

Needs to pass place & route
and pass timing closure

Operator Performance AutoTuning



Schedule Exploration with VTA

VTA Candidate Designs

#1 Design AAA @ 307GOPs

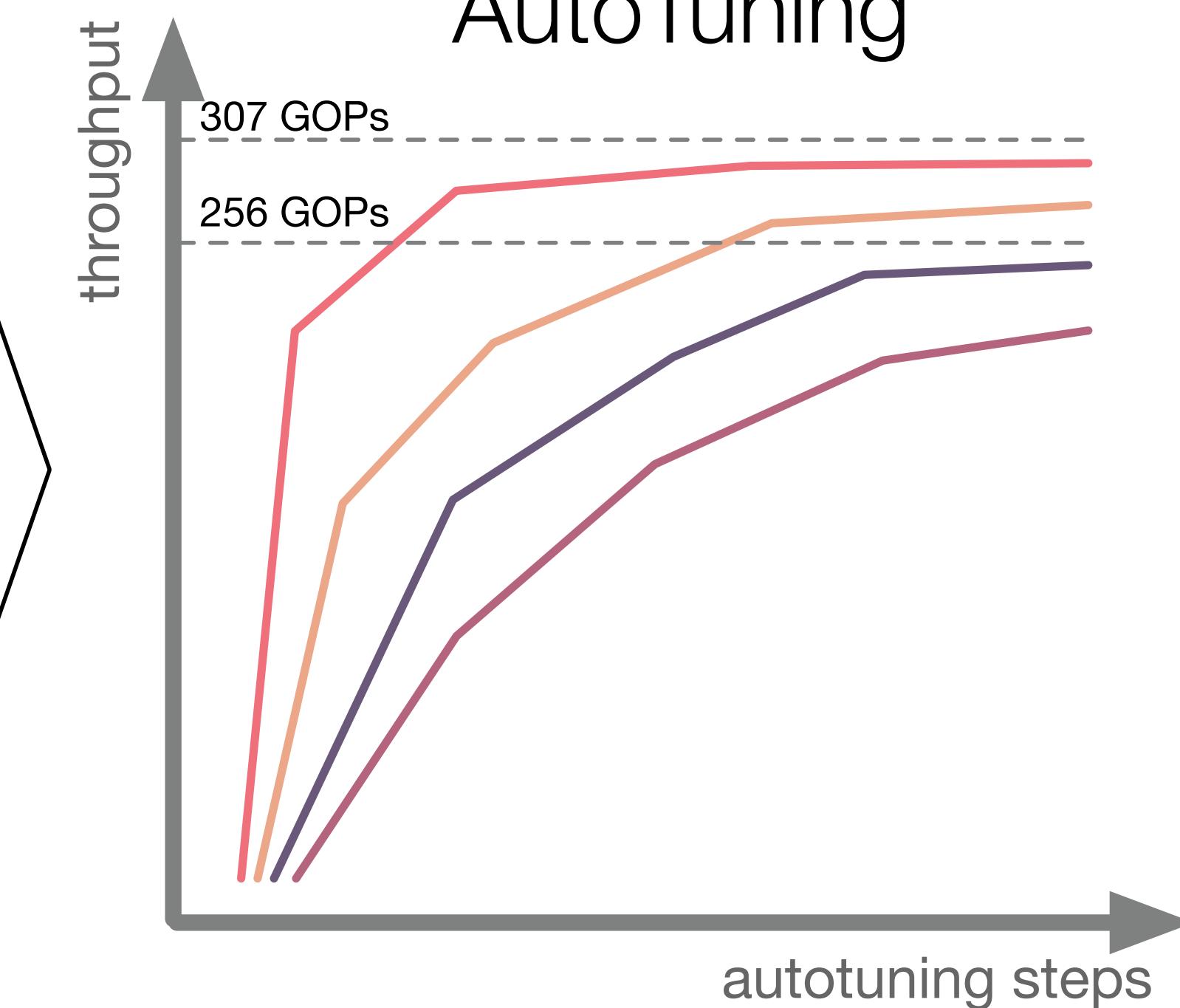
#2 Design BBB @ 307GOPs

#3 Design CCC @ 307GOPs

#4 Design DDD @ 256GOPs

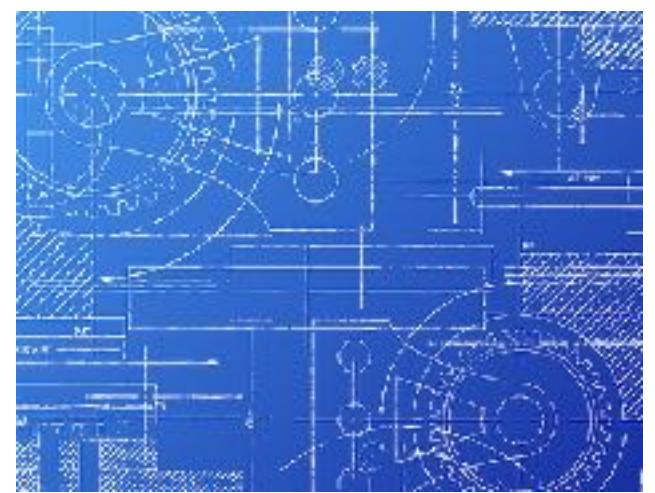
Needs to pass place & route
and pass timing closure

Operator Performance AutoTuning



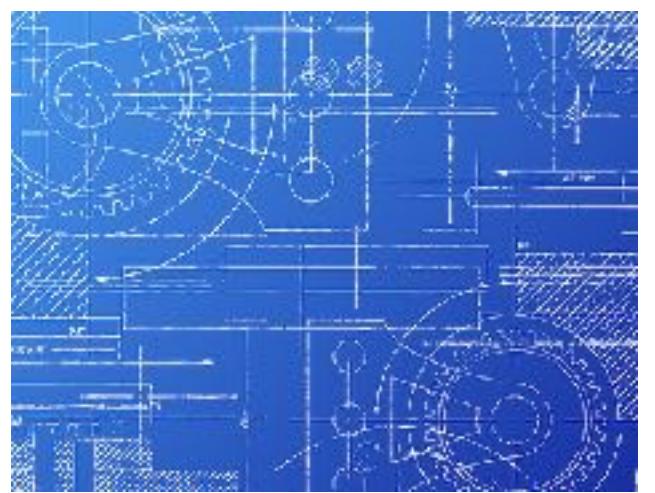
TVM+VTA Stack Goals

TVM+VTA Stack Goals



- Blue-print for a complete deep learning acceleration stack

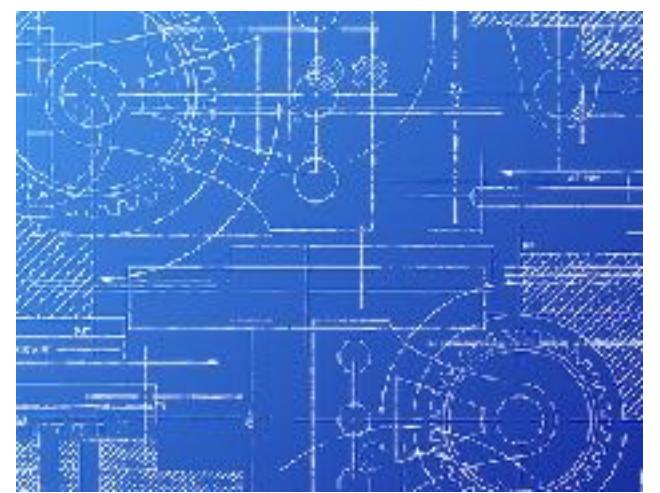
TVM+VTA Stack Goals



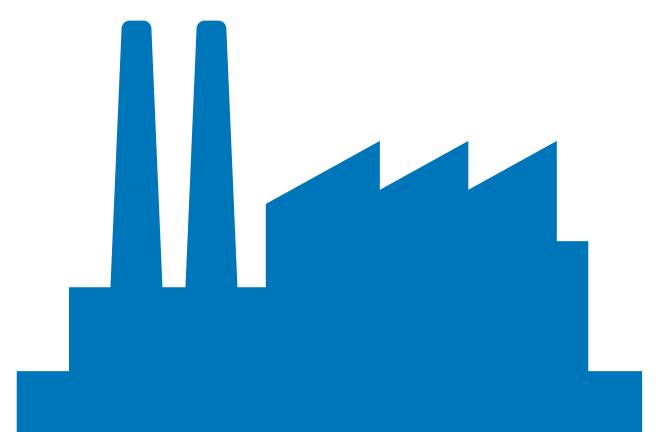
- Blue-print for a complete deep learning acceleration stack
- Experimentation framework for cross-stack deep learning optimizations



TVM+VTA Stack Goals



- Blue-print for a complete deep learning acceleration stack
- Experimentation framework for cross-stack deep learning optimizations
- Open-source community for industrial-strength deep learning acceleration

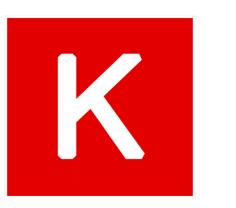
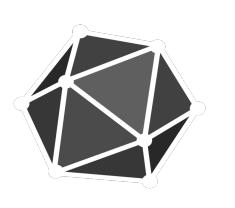
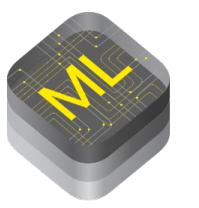


Carlos Guestrin

Training Deep Learning Models with TVM

Jared Roesch

Model



High-Level Differentiable IR

Tensor Expression IR

LLVM, CUDA, Metal

VTA



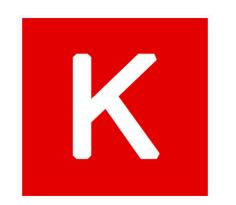
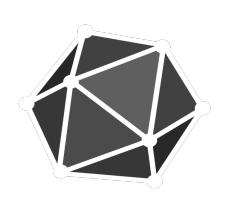
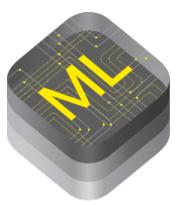
Edge
FPGA

Cloud
FPGA

ASIC

**Standalone
inference deployment**

Model



High-Level Differentiable IR

Tensor Expression IR

LLVM, CUDA, Metal

VTA



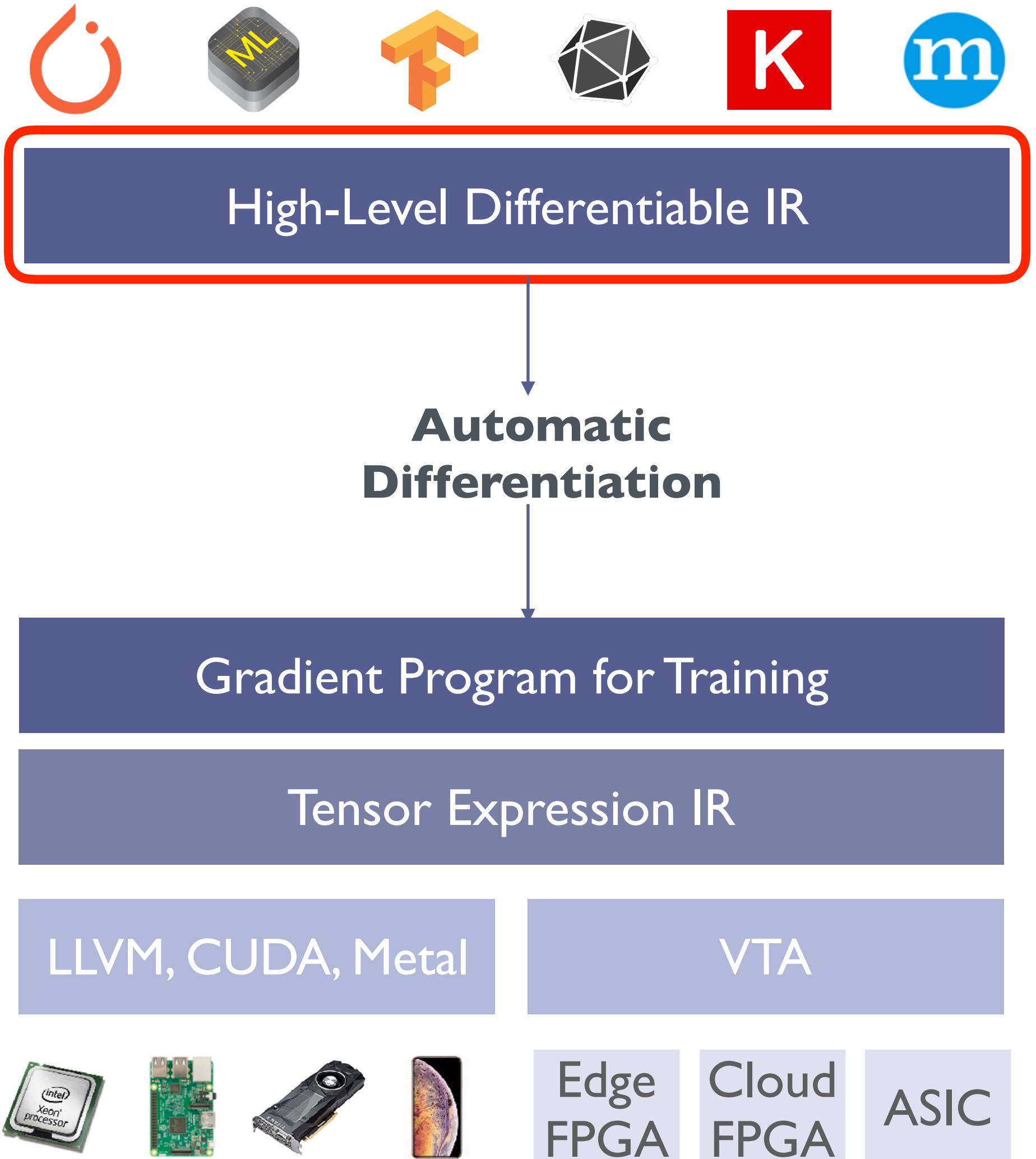
Edge
FPGA

Cloud
FPGA

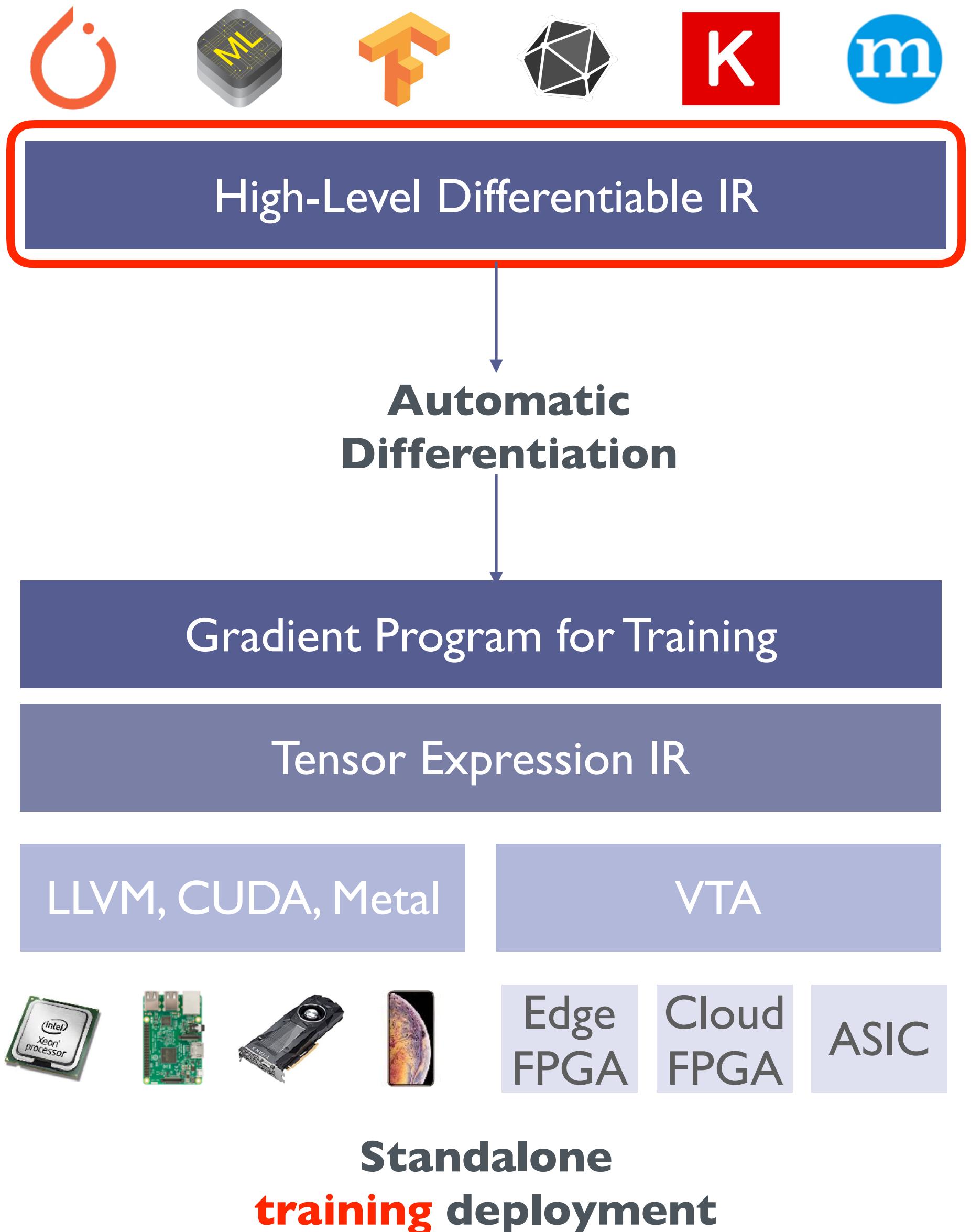
ASIC

**Standalone
inference deployment**

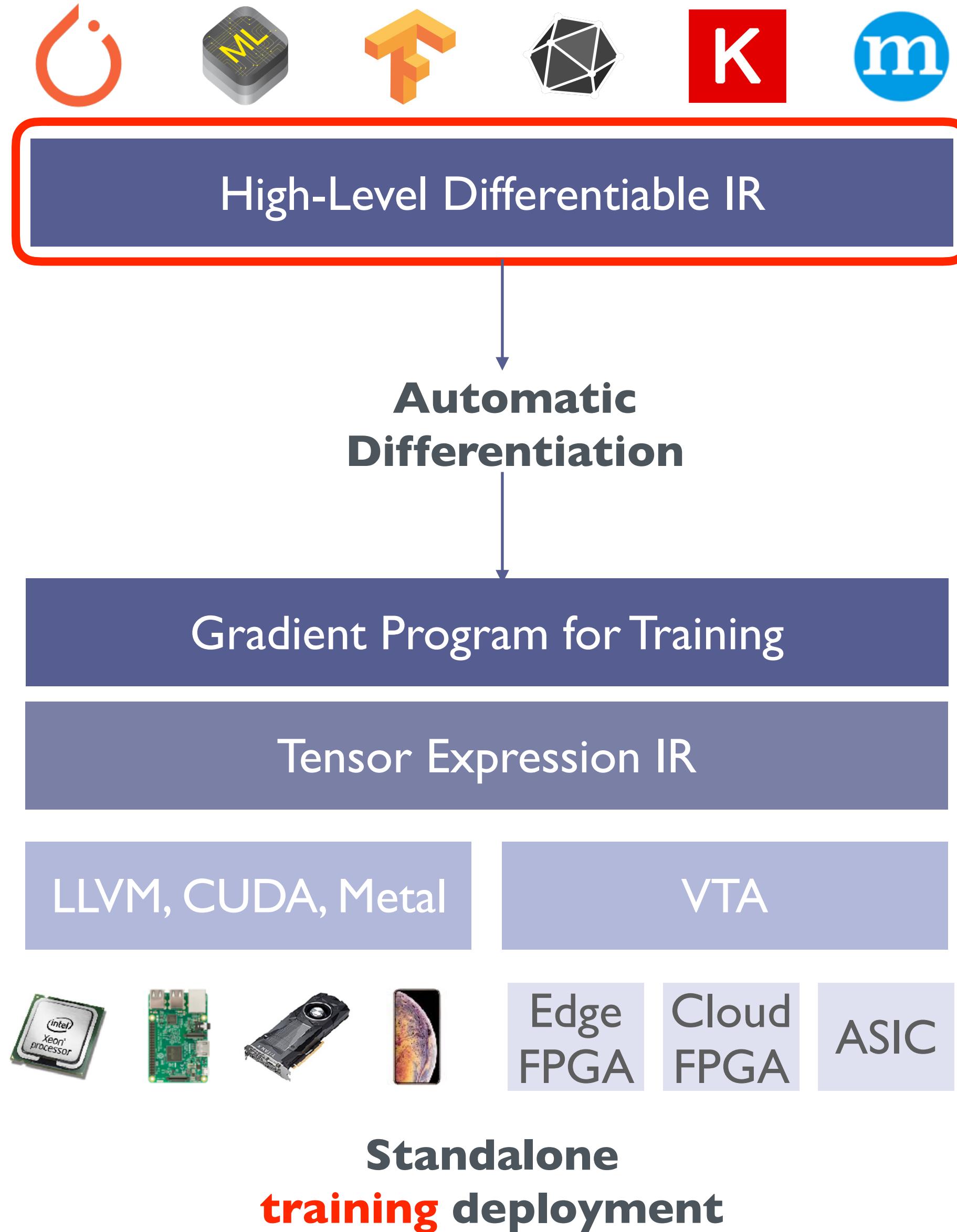
Model



Model



Model



- Automatic generation of gradient programs
- Support for customized data types and FPGA training
- Support for distributed execution, and integration with technology such as PHub (see Liang's talk).

More details on the Relay talk later today!

Road ahead...

On the horizon...

Training

Automation

Hardware

On the horizon...

Training

AutoDiff
with Relay

Training
on-device

Tradeoff accuracy/
throughput/Joules

Automation

Hardware

On the horizon...

Training

AutoDiff
with Relay

Training
on-device

Tradeoff accuracy/
throughput/Joules

Automation

Auto
quantization

Full-program
optimization

Automated
HW design

Hardware

On the horizon...

Training

AutoDiff
with Relay

Training
on-device

Tradeoff accuracy/
throughput/Joules

Automation

Auto
quantization

Full-program
optimization

Automated
HW design

Hardware

VTA Chisel
design

ASIC
flow

Training on
VTA



Big THANKS to our sponsors!



9:00	Keynote, TVM Overview, TVM @ Amazon	Keynote (SAMPL, Apple, Amazon, Huawei) TVM Overview – Tianqi Chen, UW Deep Learning Compilation at Amazon – Yida Wang, Amazon
11:05	<i>Break</i>	
11:25	Automation, HW Specialization, Security	AutoTVM & Device Fleet – Eddie Yan, UW VTA Open Source Deep Learning Accelerator – Thierry Moreau, UW Secure Enclaves for Deep Learning – Nick Hynes, UC Berkeley/Oasis Labs
12:30	<i>Boxed lunches</i>	
13:30	Training, Programming Systems, Hardware	Kunle Olukotun/Raghuram Raghu Prabhakar, Stanford & SambaNova Machine Programming – Justin Gottschlich, Intel PlaidML Stripe: Polyhedral IR + Model-guided Optimization – Brian Retford, Intel The Relay Differentiable IR for TVM – Jared Roesch, UW Scalable Distributed Training with Parameter Hub – Liang Luo, UW The HammerBlade ML Supercomputer – Michael Taylor, UW
15:20	<i>Break, contributors meetup</i>	
15:50	Compilers, FPGAs	Andrew Tulloch, Facebook Graham Schelle, Xilinx
16:30	Lightning talks	
17:35	Community formation	Markus Weimer, Microsoft and Apache Software Foundation
18:10	Social (food, drinks)	
20:00	<i>adjourn</i>	