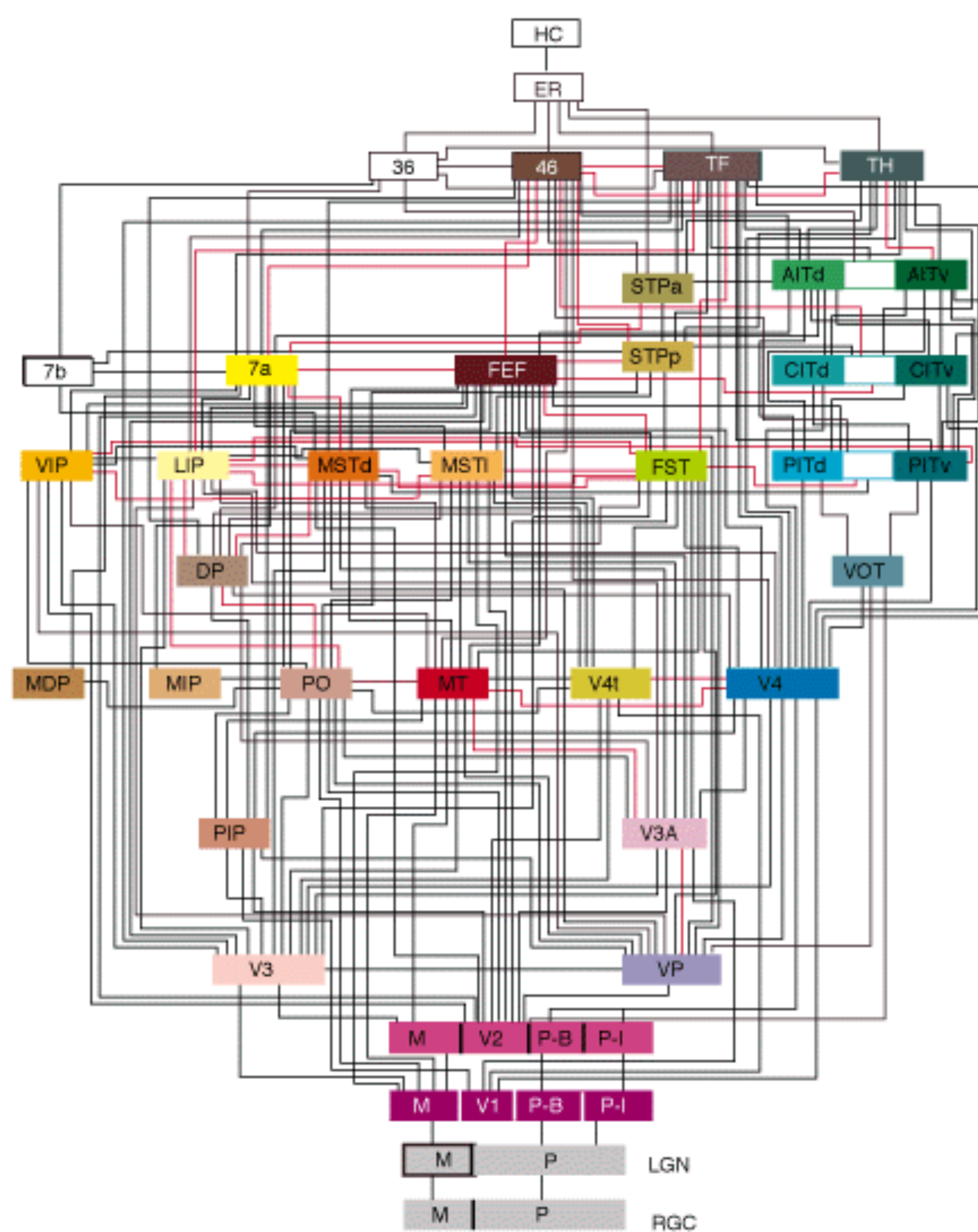


# Recent Trends in Computer Vision and Deep Learning Systems

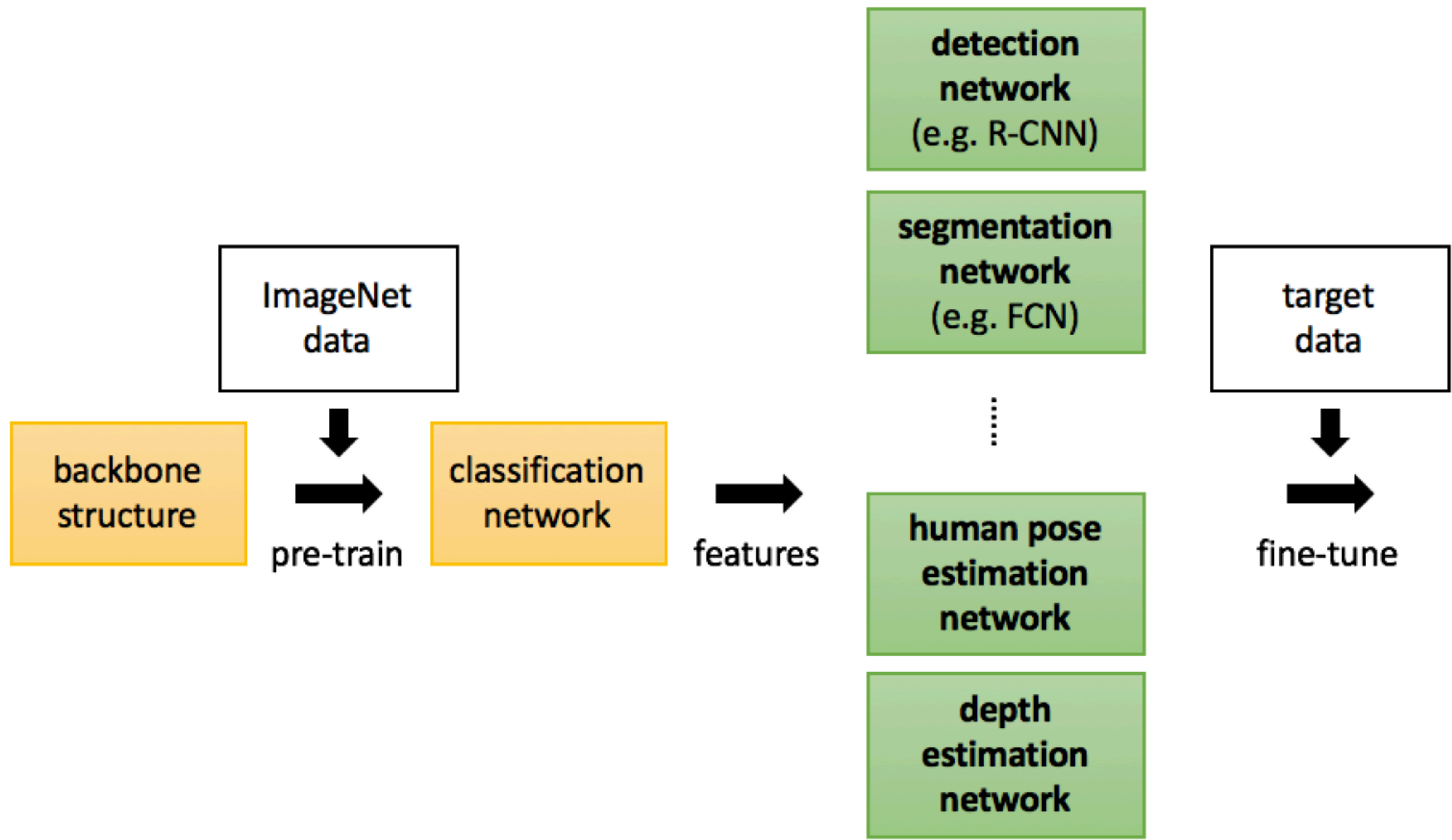
**Yangqing Jia**

Lead Researcher and Manager of AI Platform, Facebook



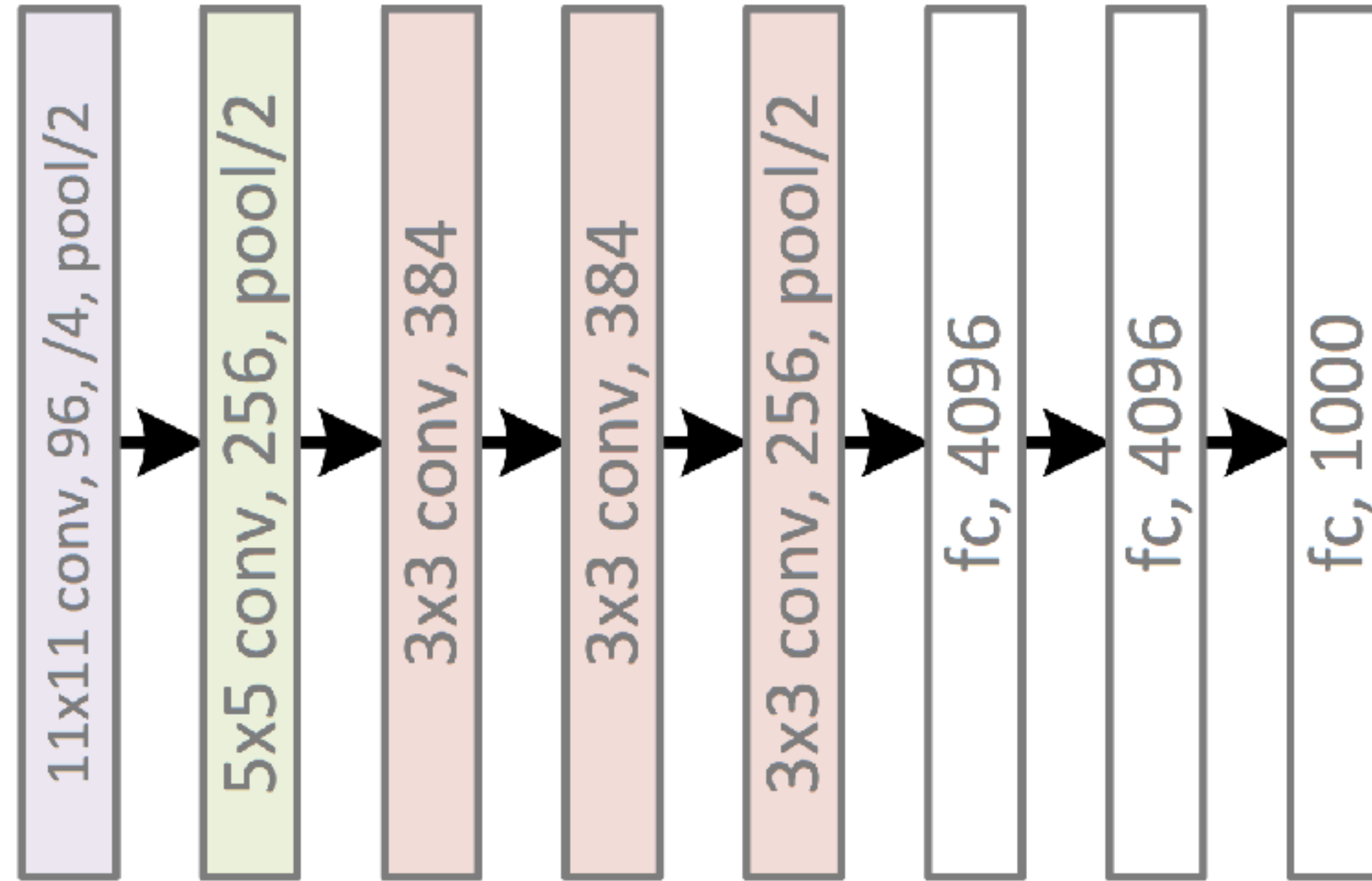


# Computer Vision



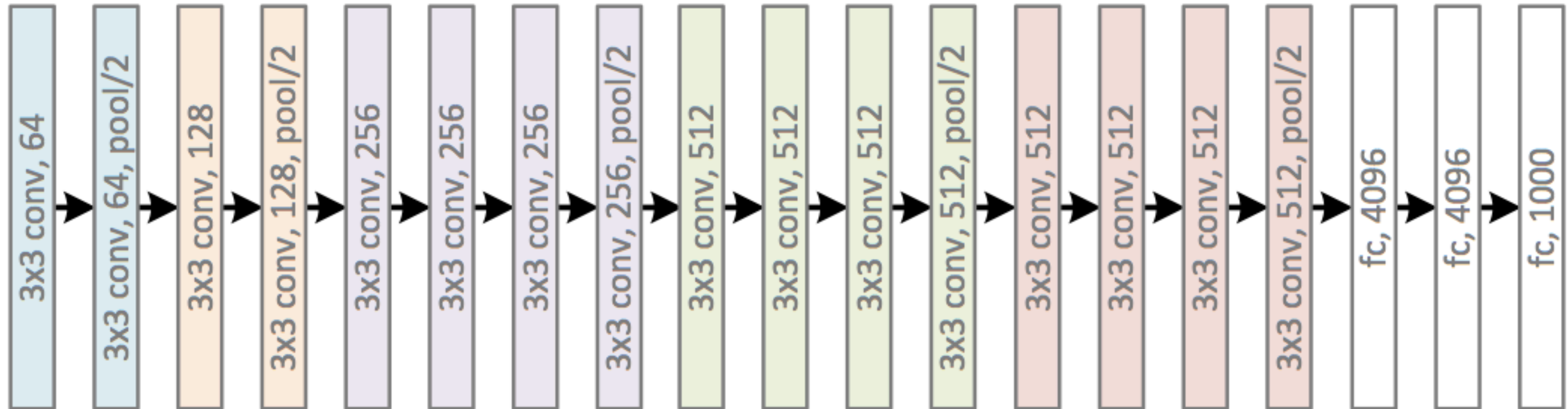
# AlexNet

So it begins.



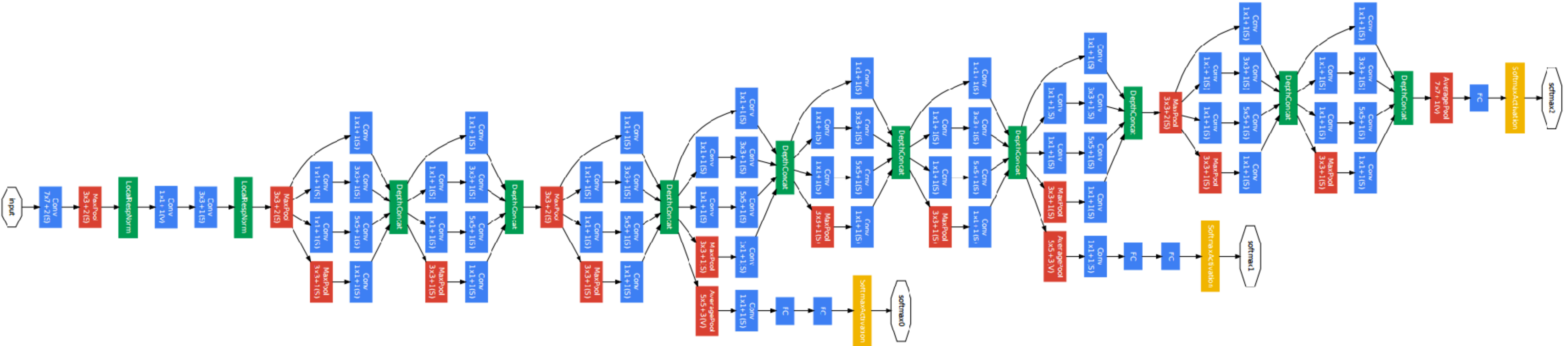
# VGGNet

Punch it.



# GoogLeNet

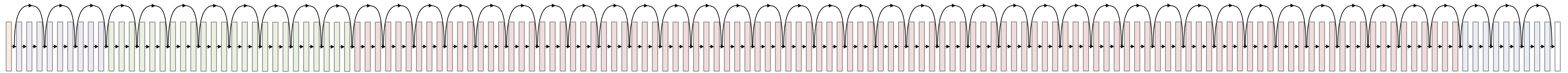
We must go deeper.





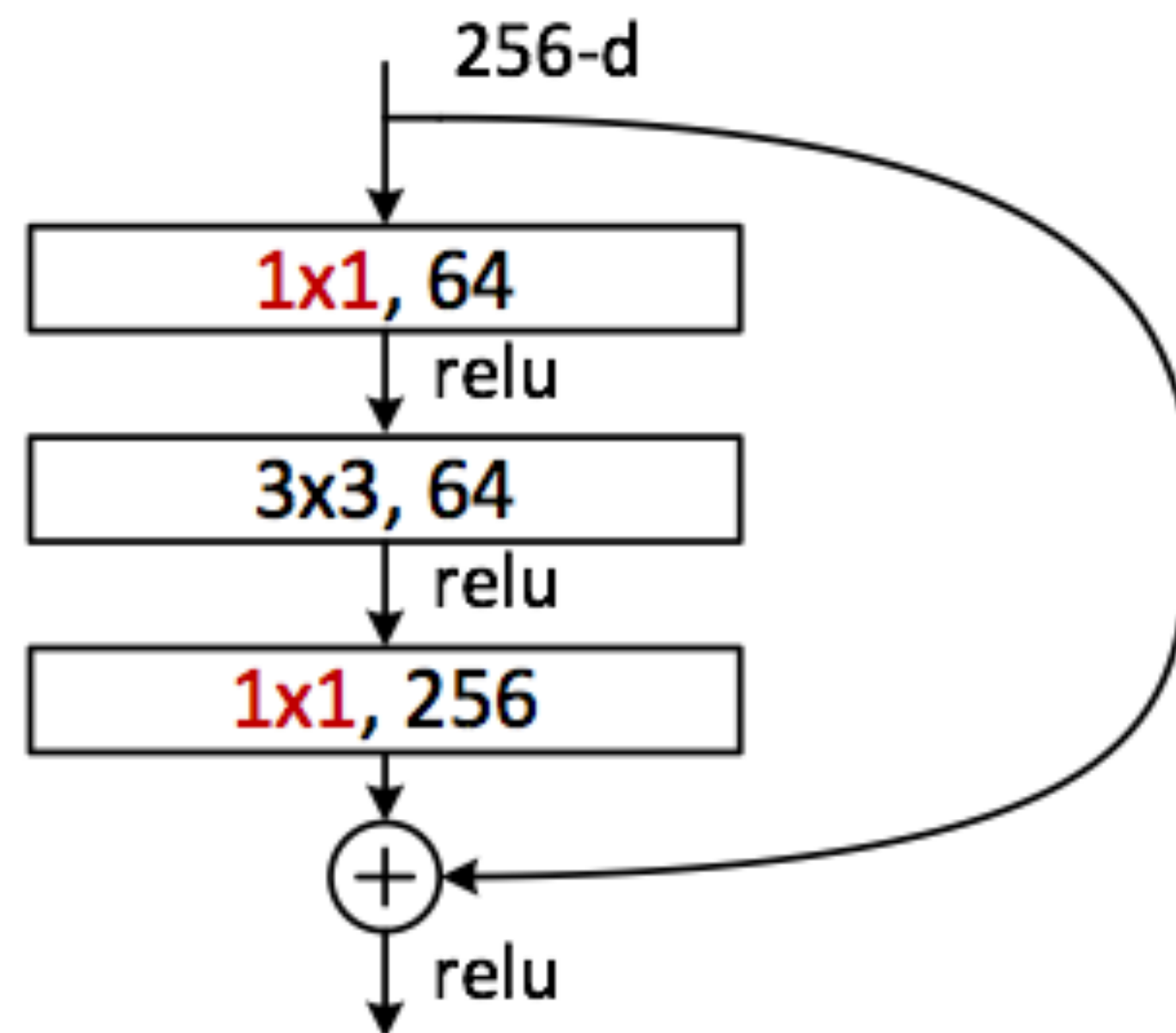
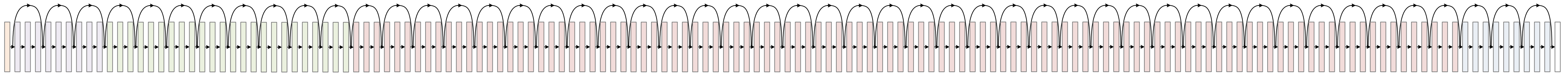
# ResNet

And we took the word seriously



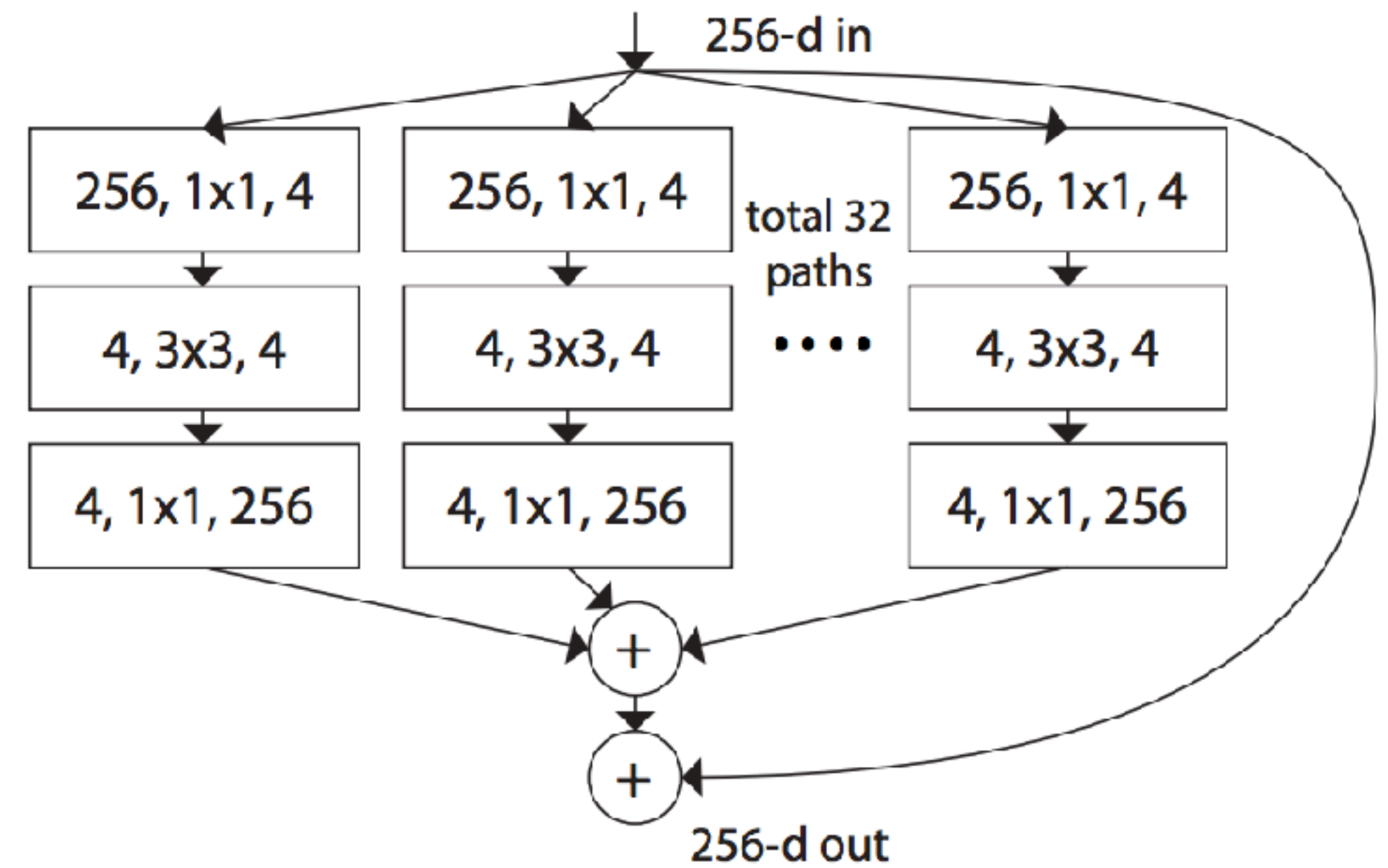
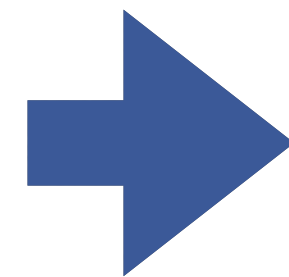
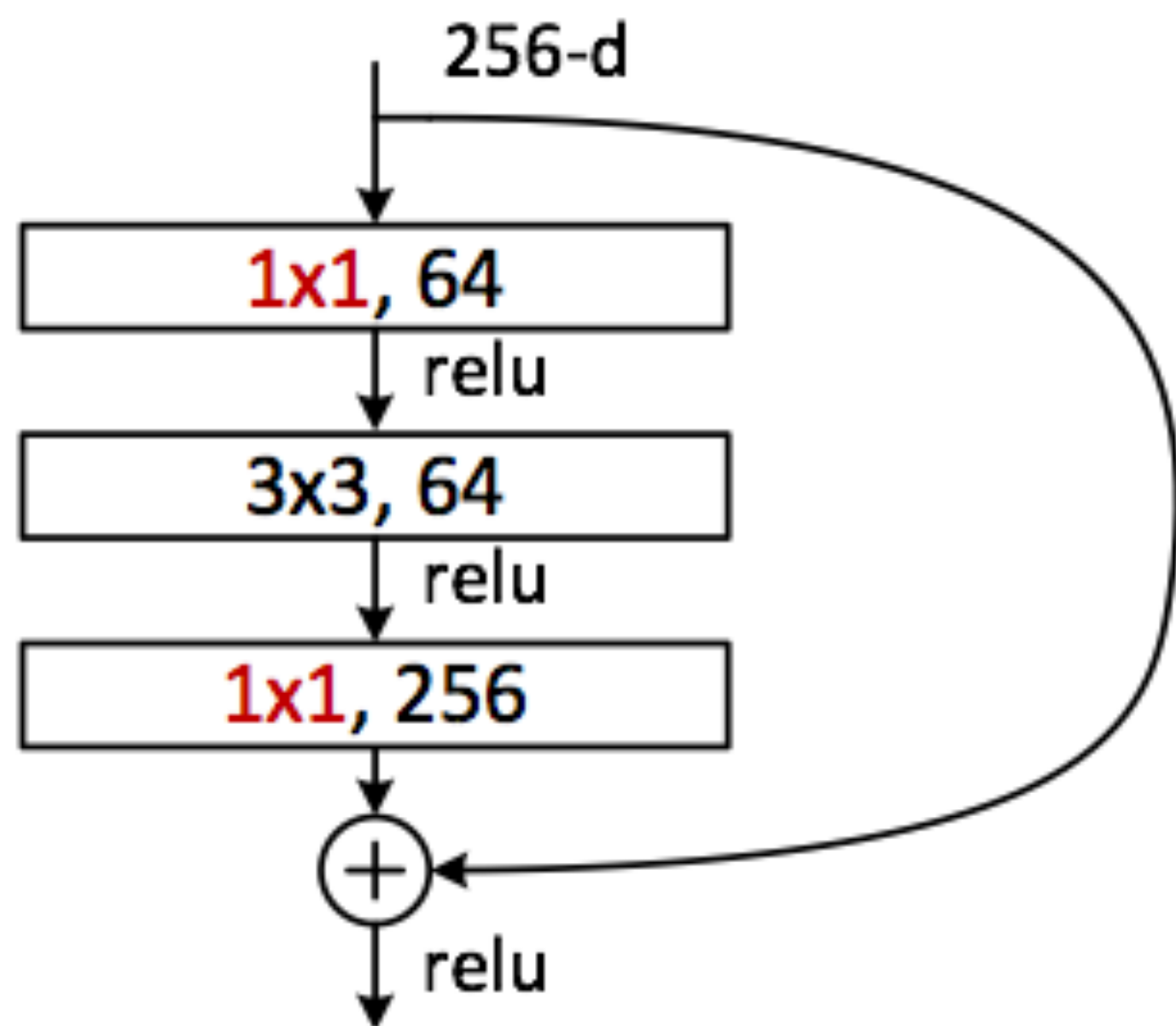
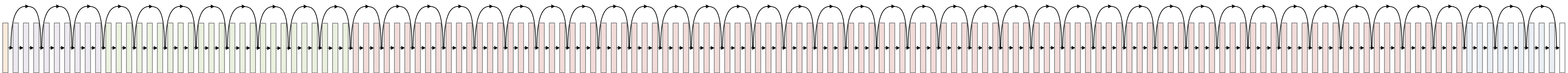
# ResNet

And we took the word seriously

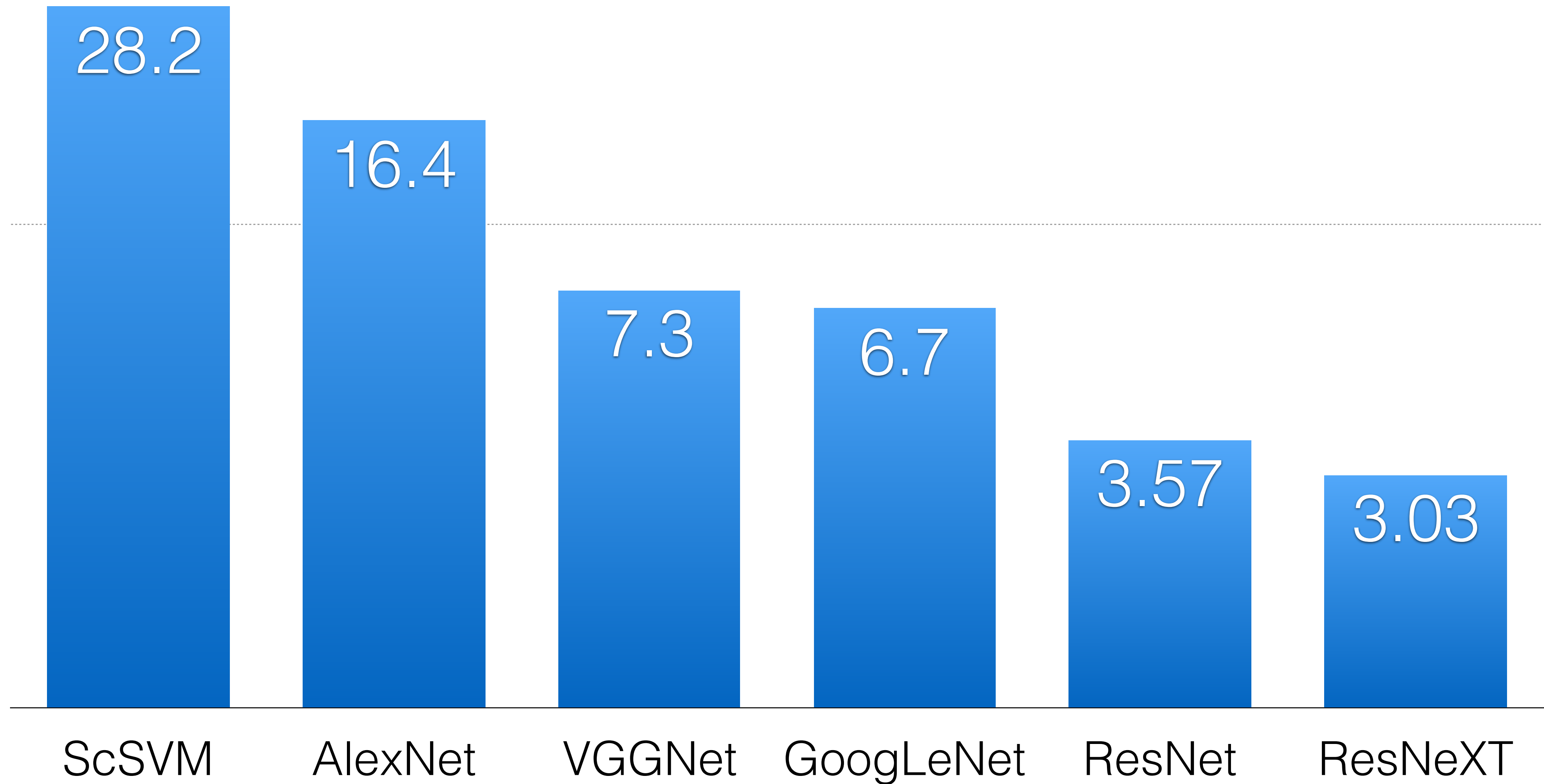


# ResNeXT

We totally see it coming

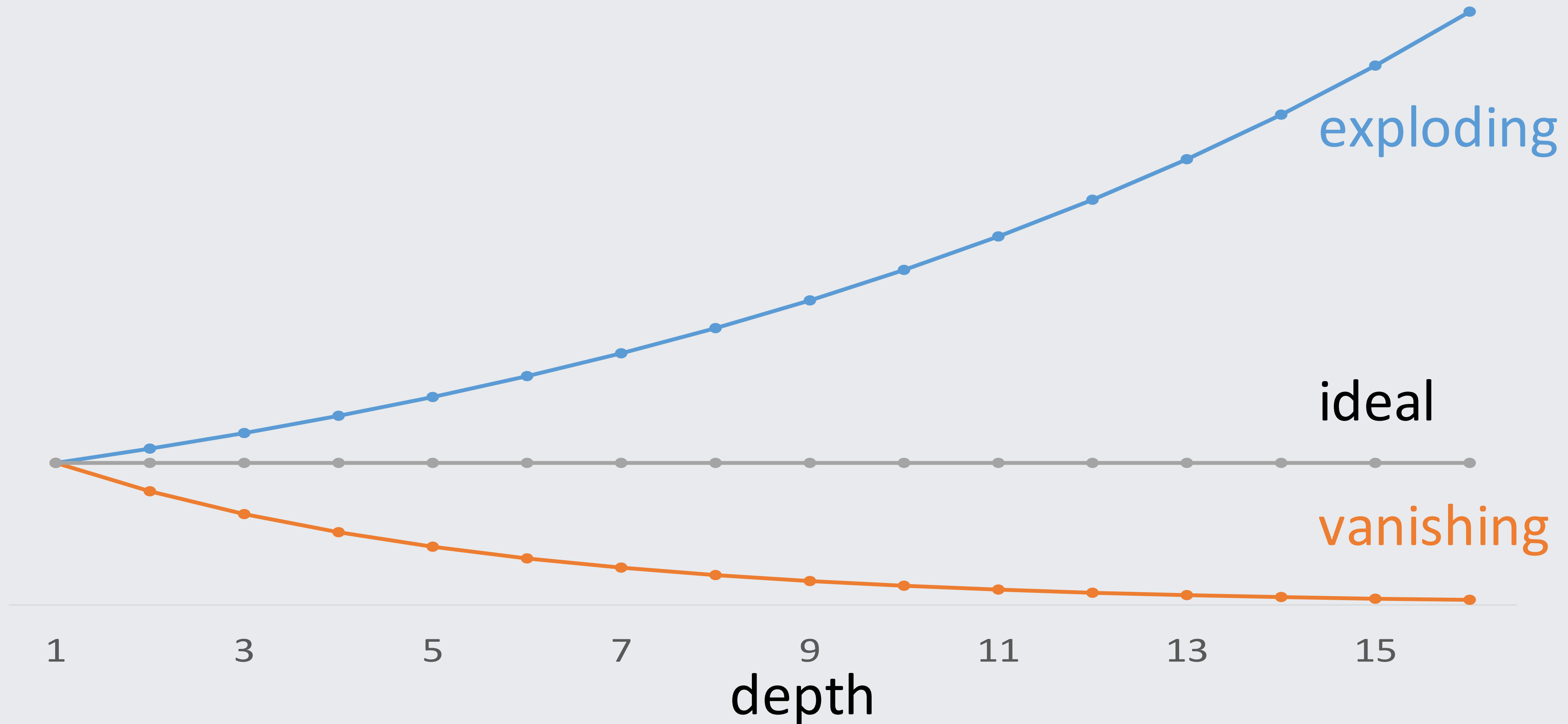


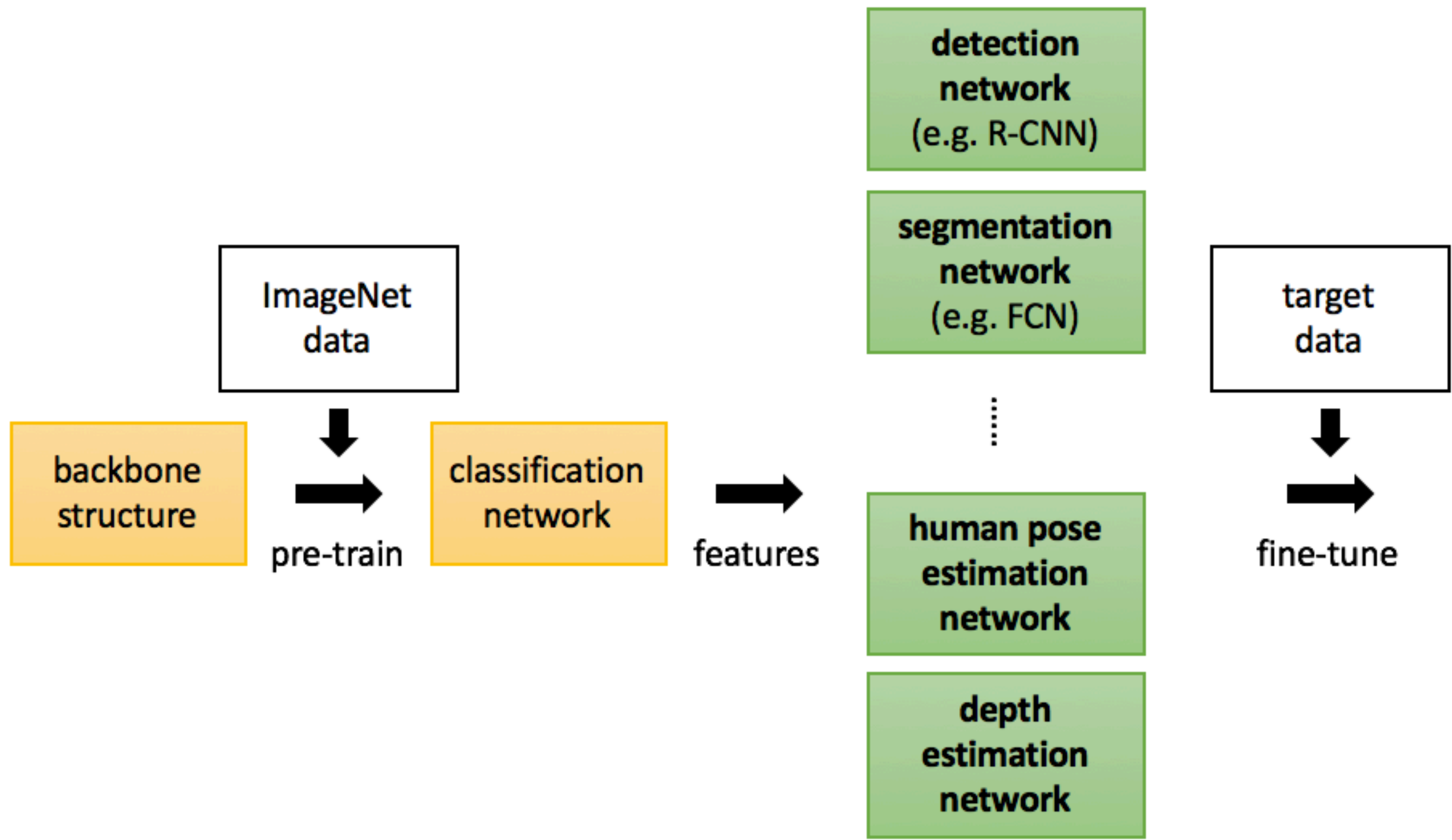
# Pushing the Performance



# Why is it challenging?

Gradients, as one example











# Deep Learning Systems

“SAP”

-

Scalability

# Scalability

Run fast, run far

“How do I train on multiple GPUs and machines?”

- Probably the most question we got from Caffe users

# Scalability

Run fast, run far

1.2 million =

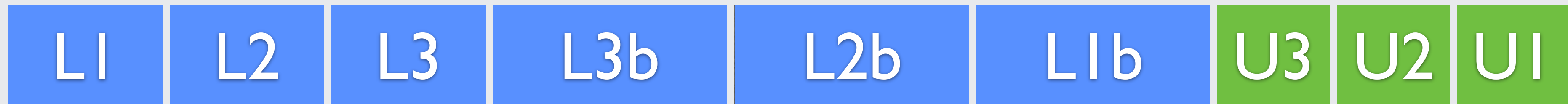
(# of images in ImageNet1K)

(# of new images @FB every 5 mins in 2013)

(# of AI jobs per month @FB)

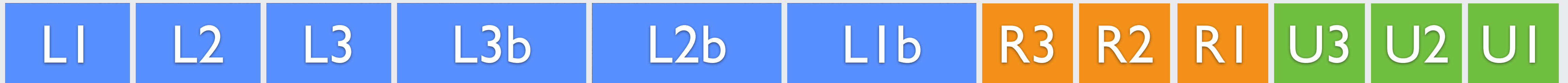
# Scalability

Run fast, run far



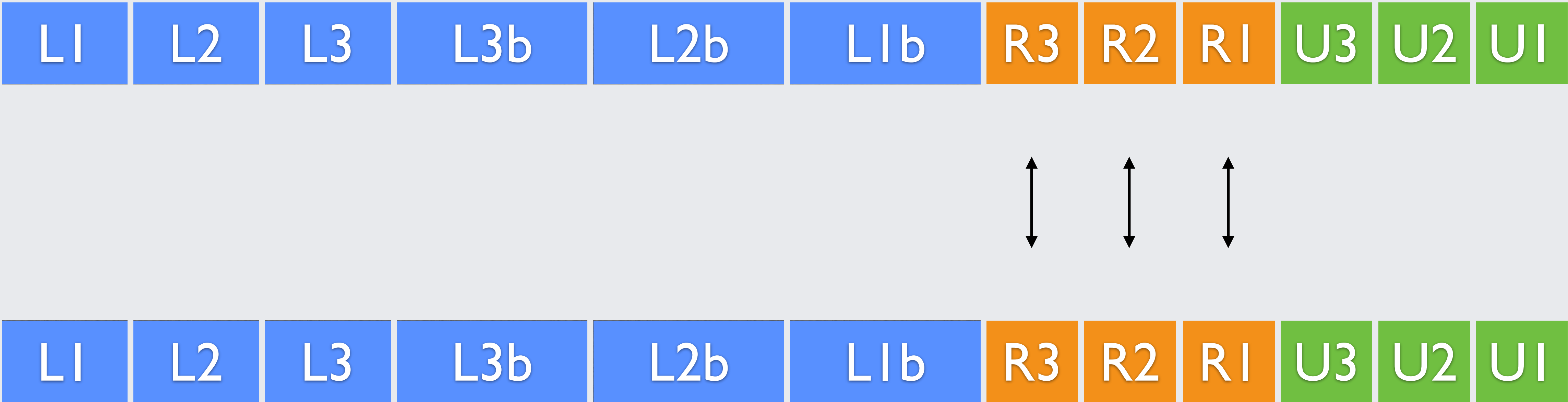
# Scalability

Run fast, run far



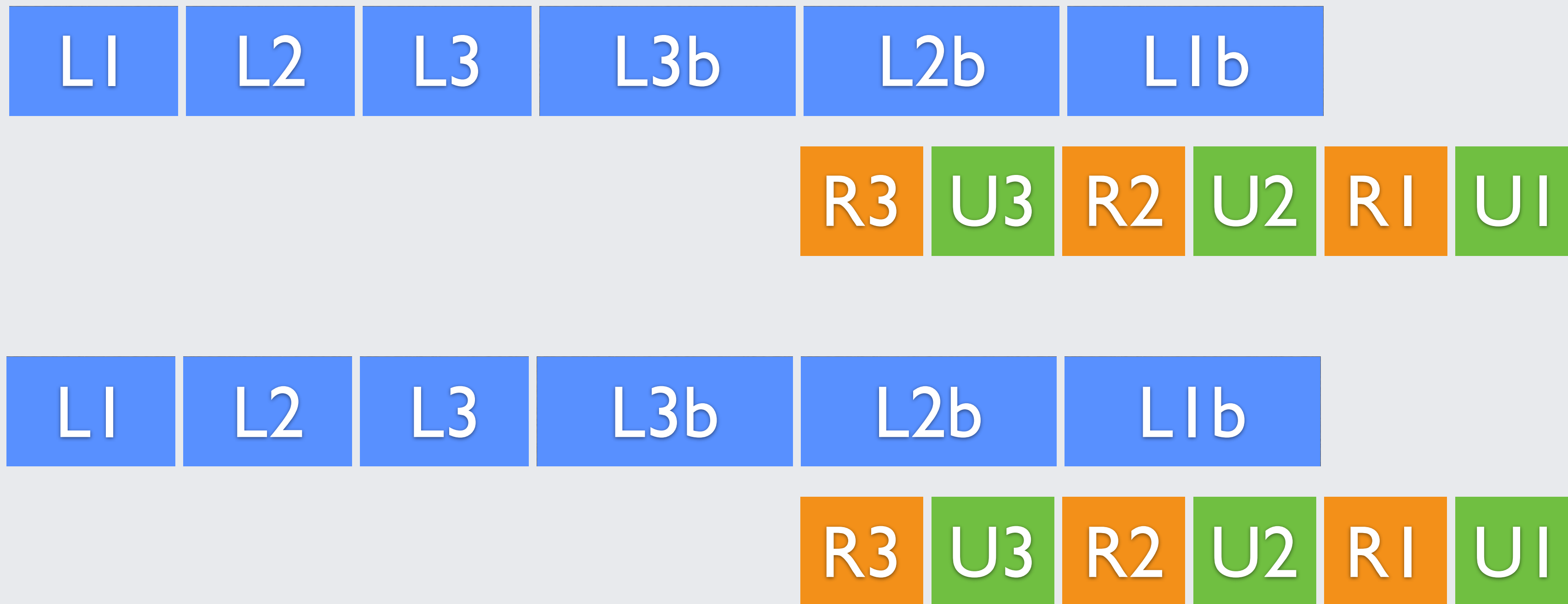
# Scalability

Run fast, run far



# Scalability

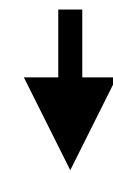
Run fast, run far





# The Return of MPI

"I'm your father", said Allreduce.

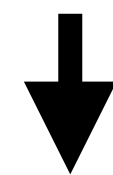


Allreduce

Tree based -  $O(M \log N)$

Ring based -  $O(M)$

etc.

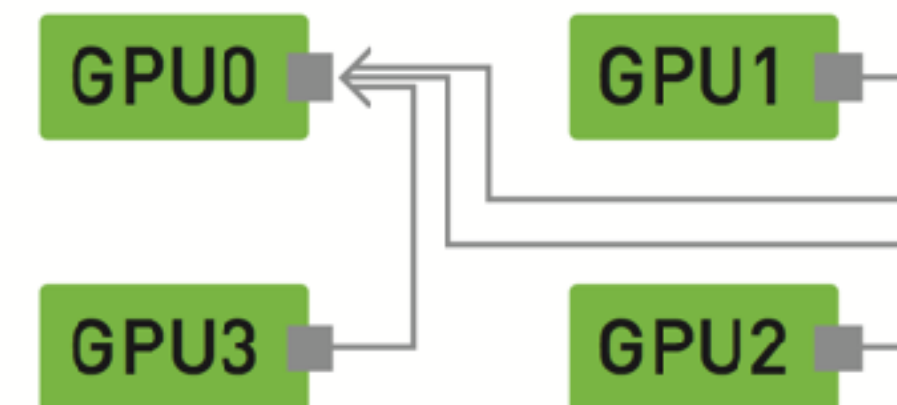


OPEN MPI

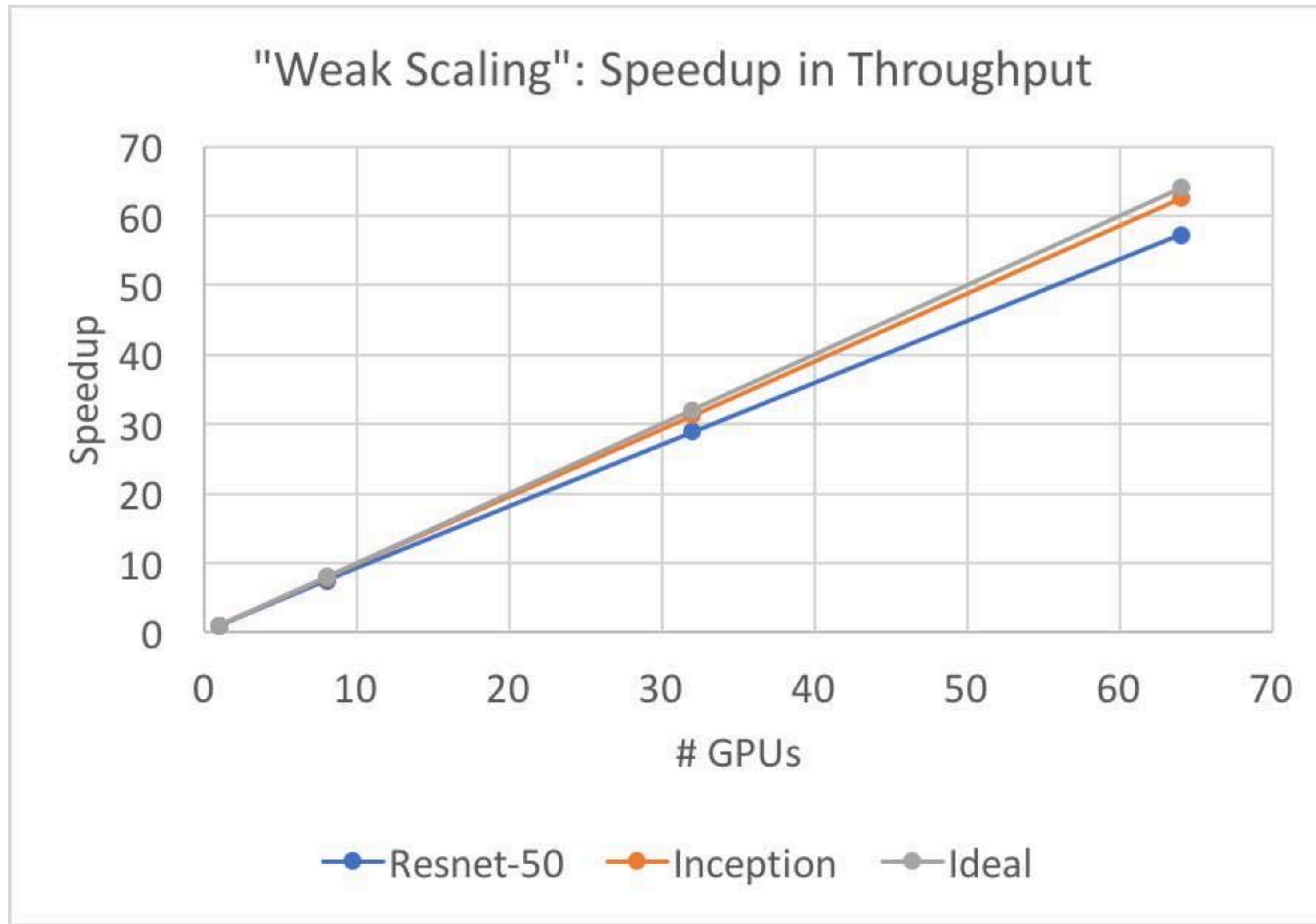


**MVAPICH**

NCCL



# And so we scale



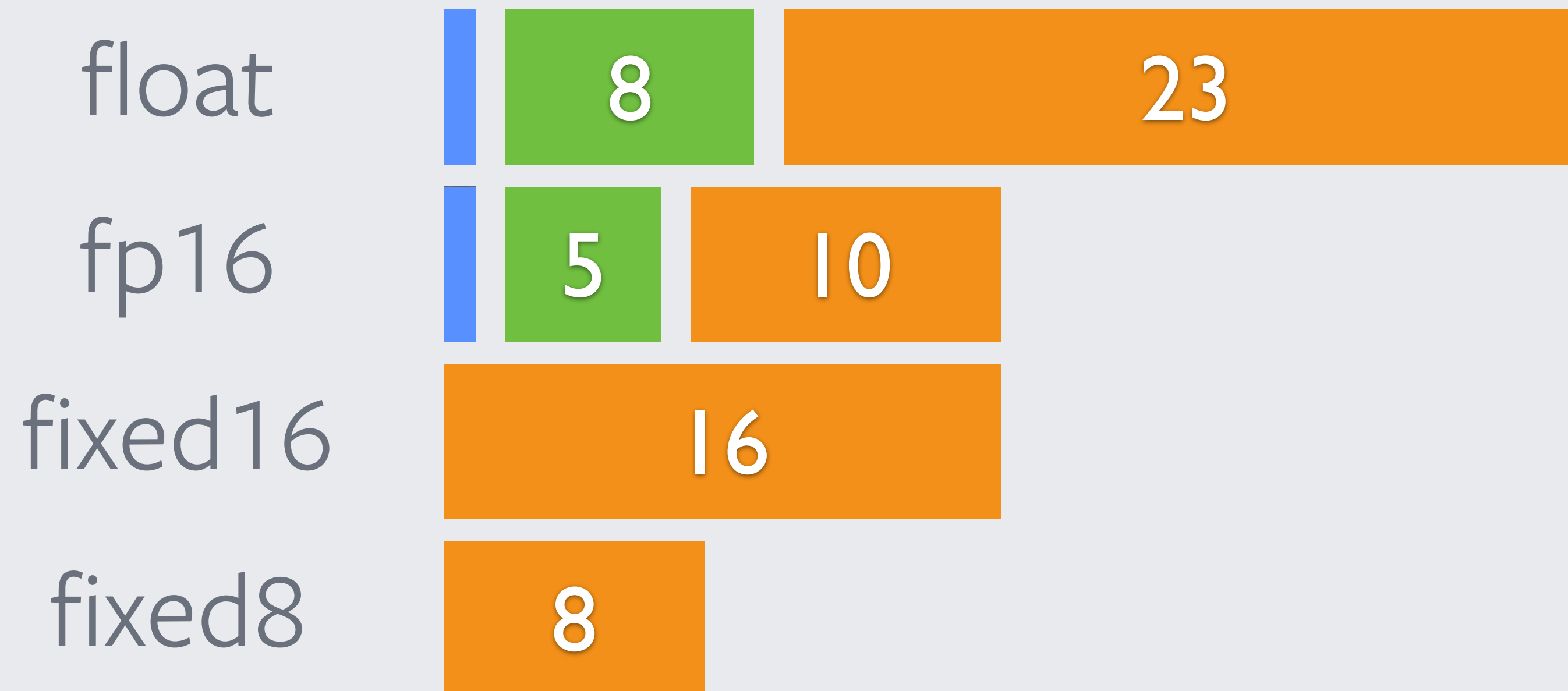
"SAP"

-

Arithmetics

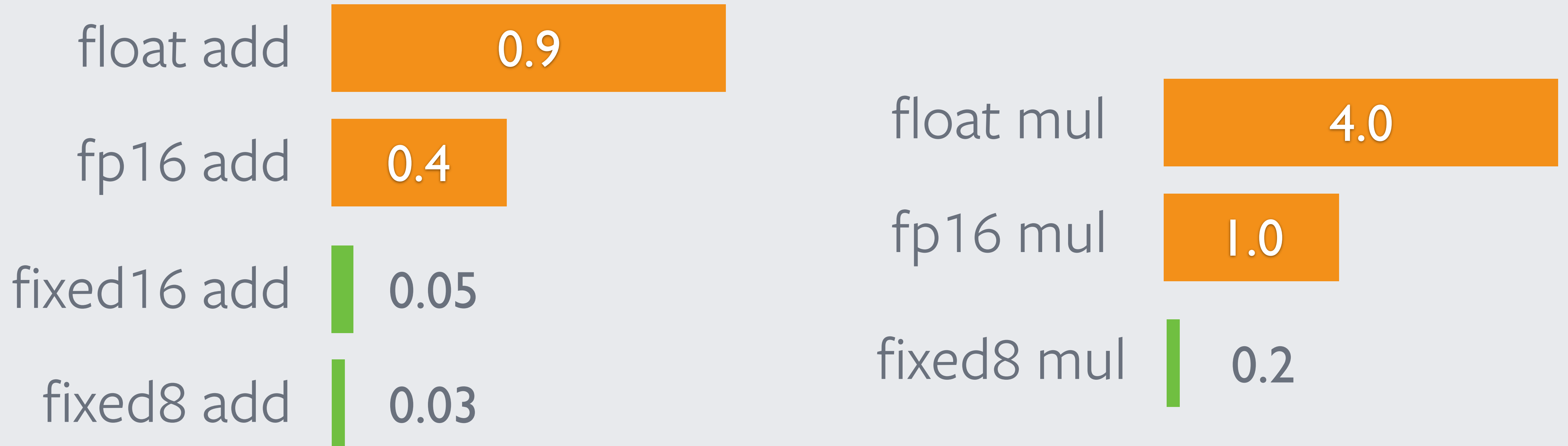
# Quantized Computation

Forget about float, the world is bigger

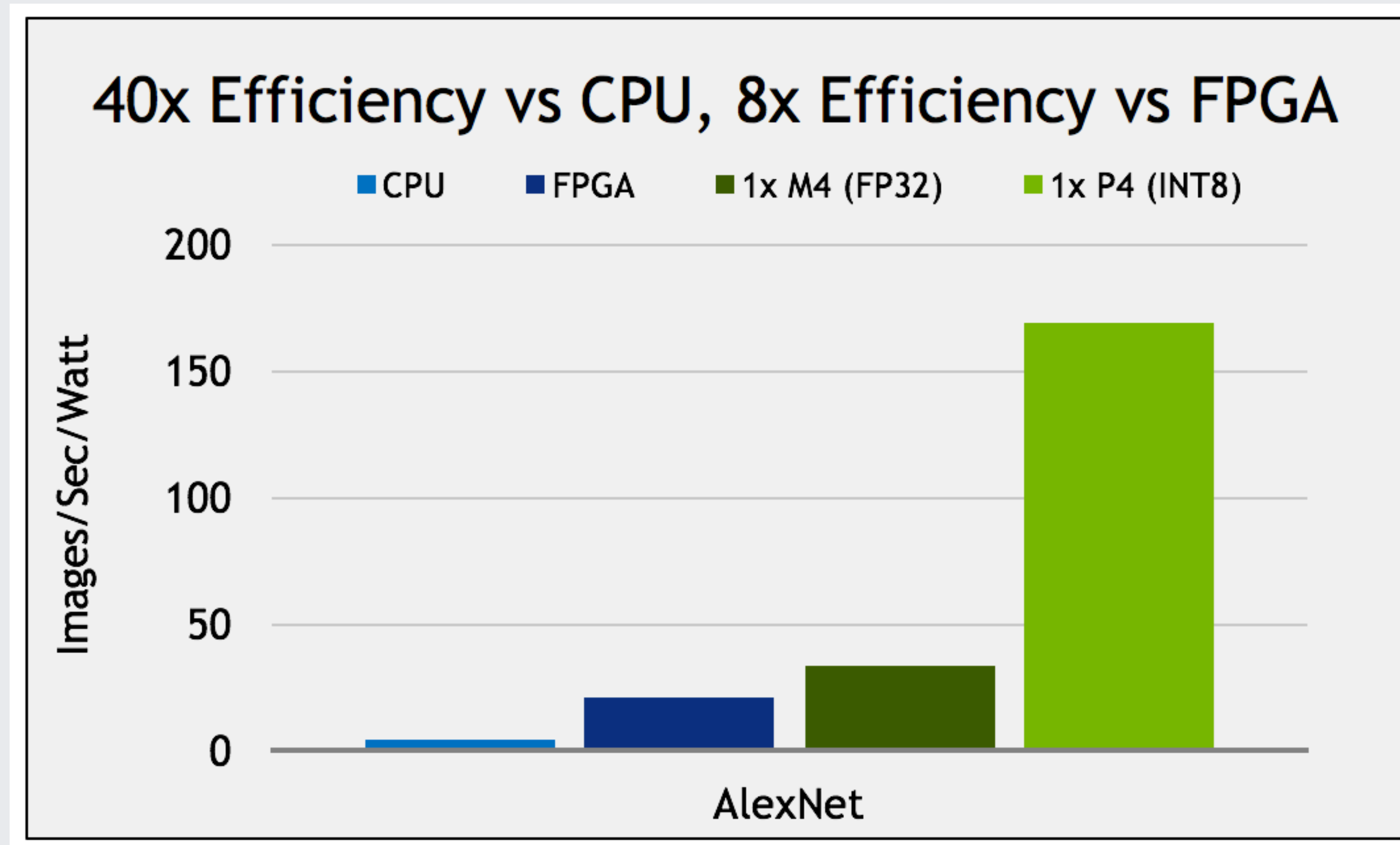


# Why do we care?

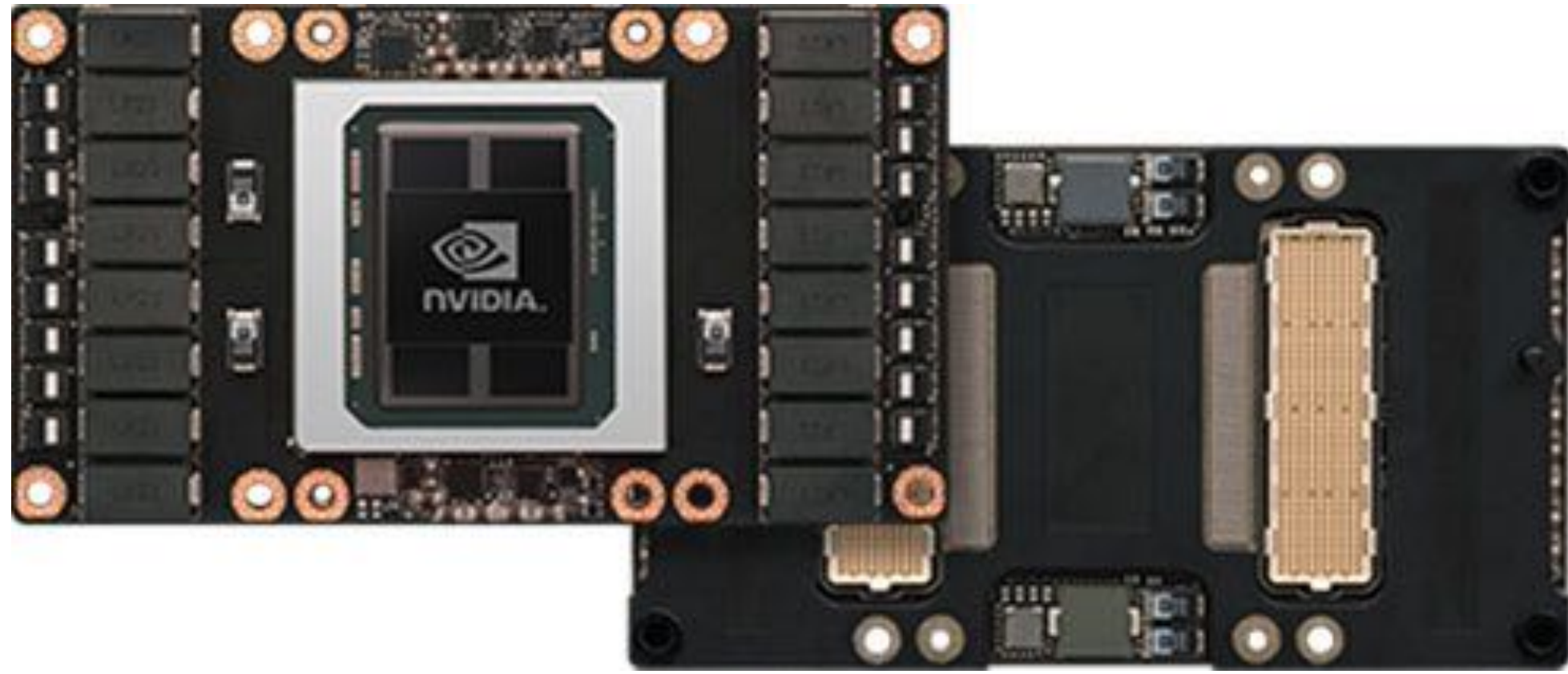
Battery life is life.



# How does it perform?



# Why does it matter for cars?



**250 watts**  
**10 -> 20 TFlops**



**10 watts**  
**0.7 -> 1.5 TFlops**

**"SAP"**

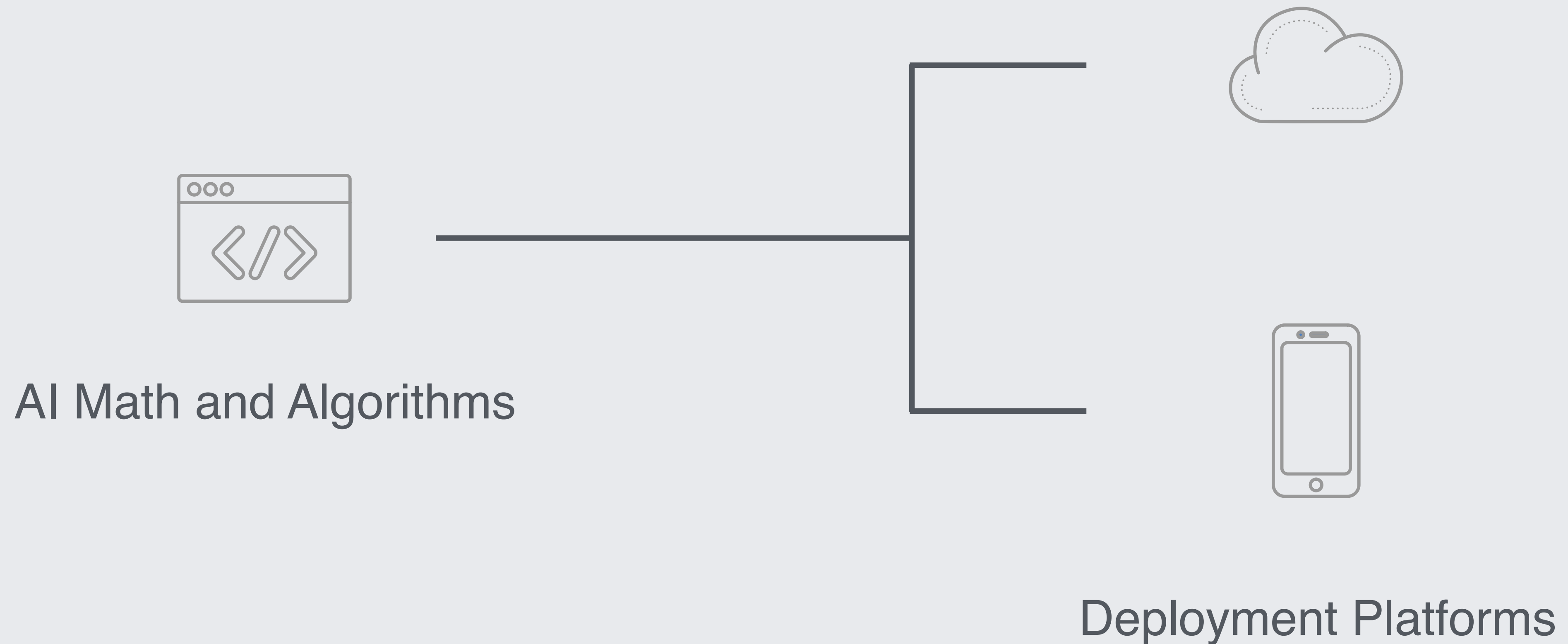
-

**Portability**



# Portable System

One software to rule them all, and...





# Portable System

Cloud, Mobile, IoT, Cars, Drones, Coffee makers





# The Land of Deep Learning System

Not as complex as a car, but still.

Applications

Caffe, Torch, TF, etc...

DataBases

LevelDB  
RocksDB  
Hadoop  
Amazon S3  
your old disk

Core Math

Eigen  
CuDNN  
NNPack  
THNN  
MKL

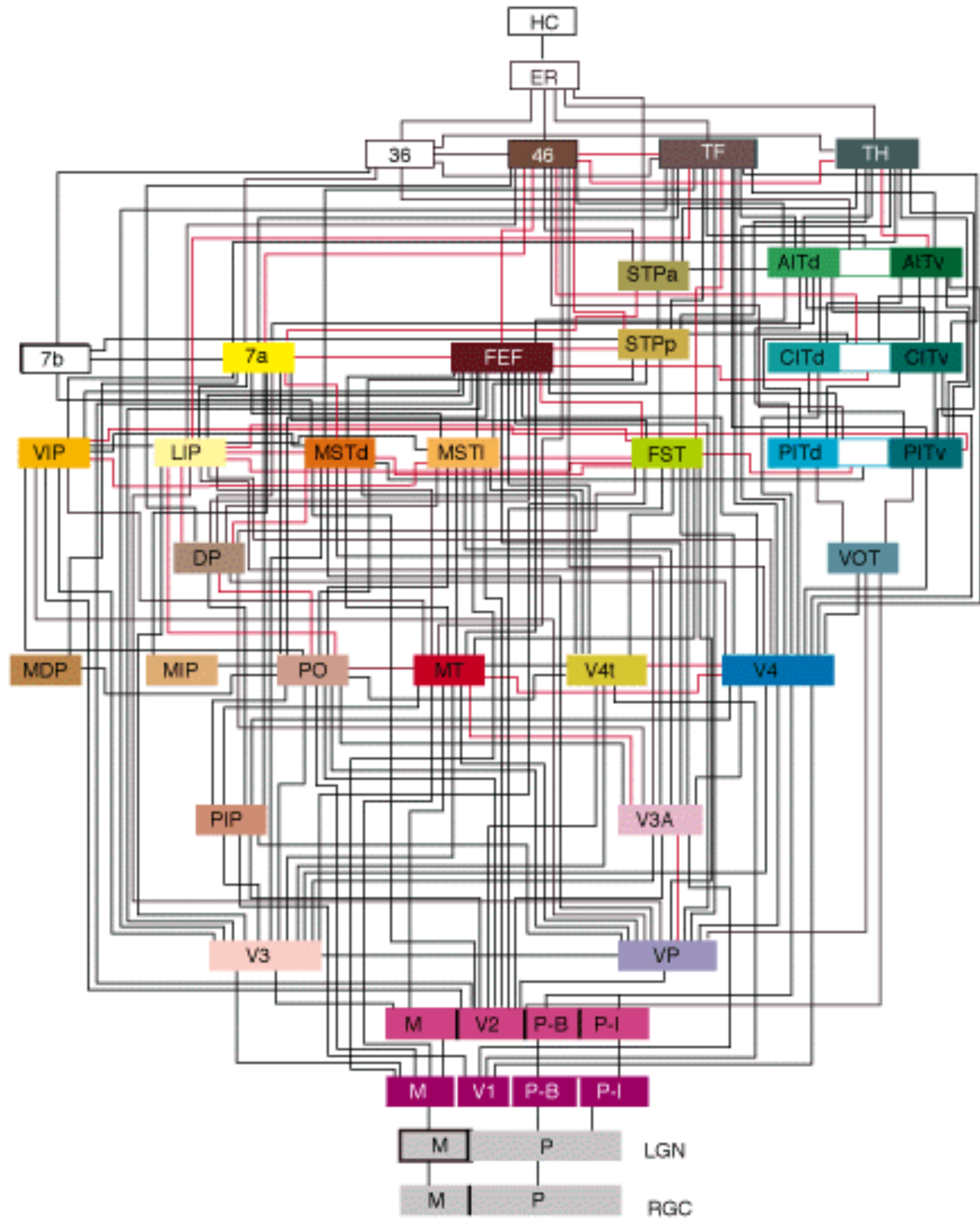
Comms

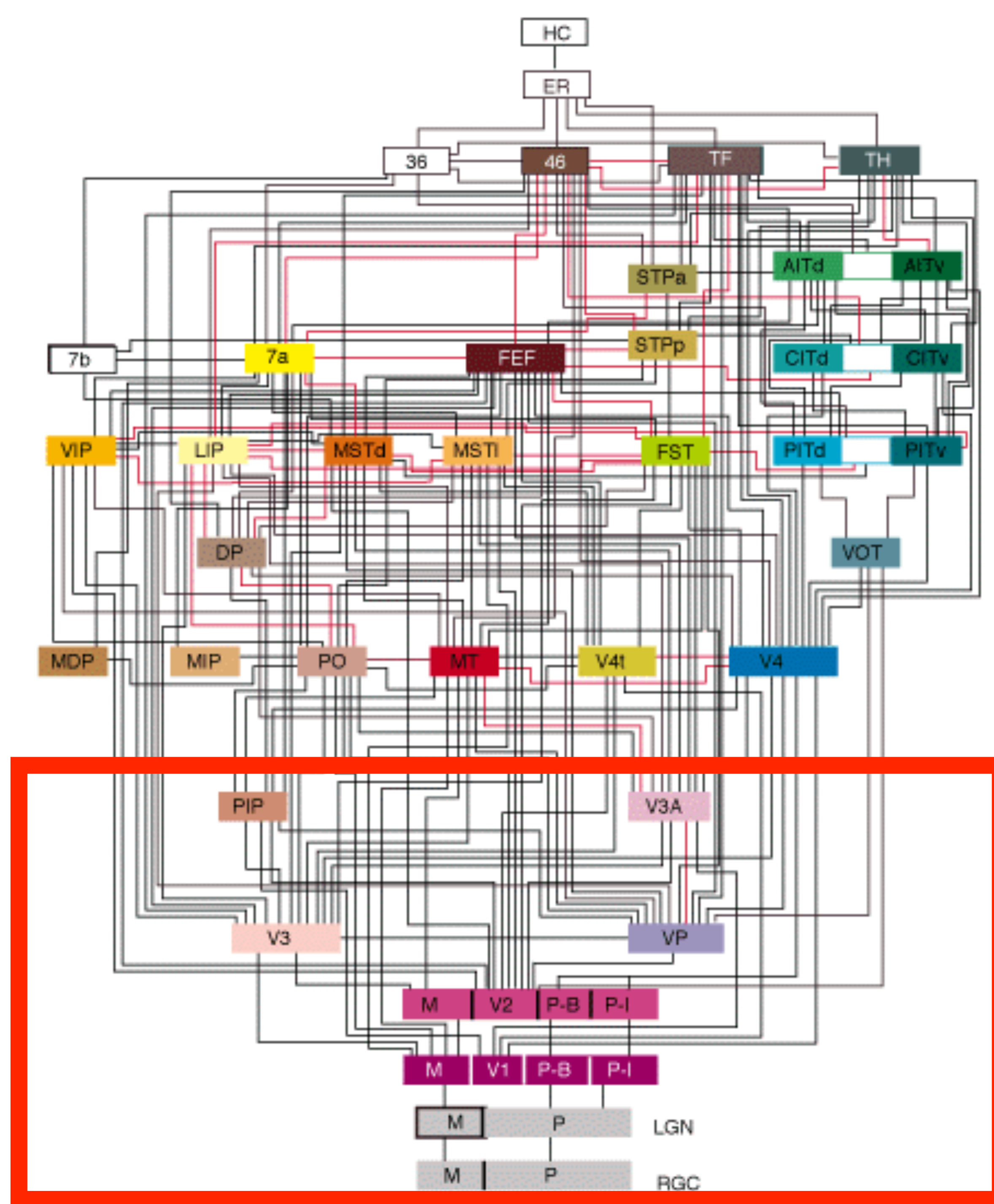
NCCL  
MPI  
ZeroMQ  
Redis  
...

Low Level

CUDA  
OpenGL  
OpenCL  
Vulkan  
...

Compilers





# Thank you!

**Recent Trends in Computer Vision and Deep Learning Systems**

Yangqing Jia