

# **Understanding Text on Images with AI at Scale**

**Viswanath Sivakumar**

Facebook AI Research (FAIR)



PINGUINO FELIZ TE ESPERA

**PINGUINO FELIZ TE ESPERA**

EN MAGDALENA EN 2018

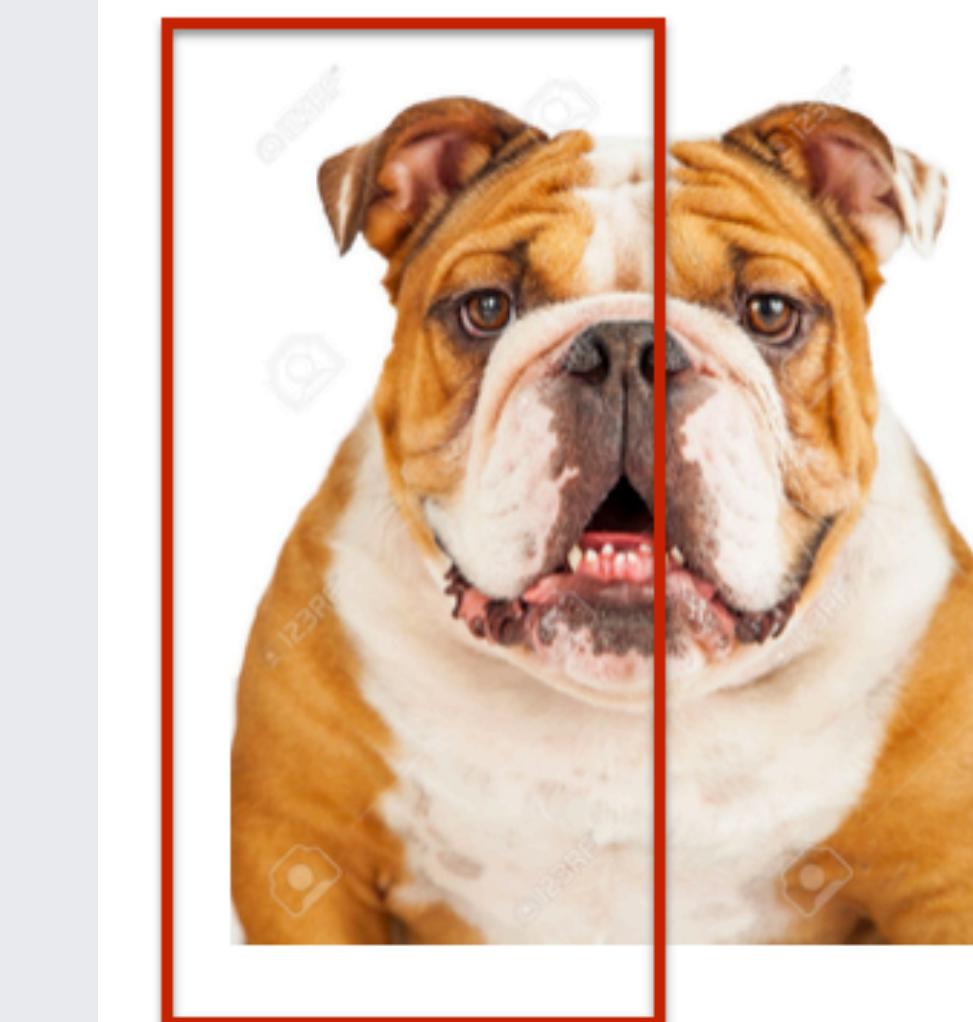
**EN MAGDALENA EN 2018**

VIENES?

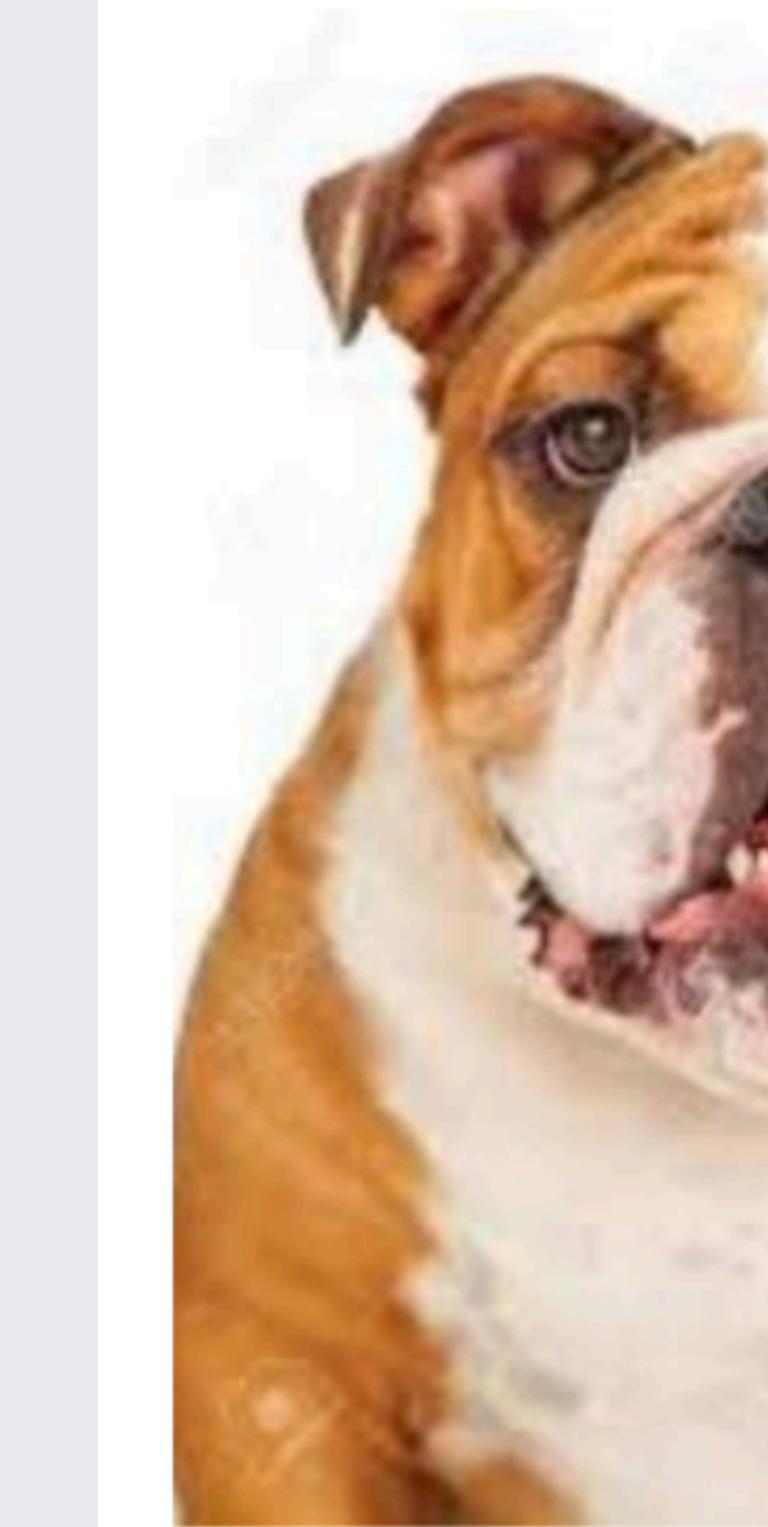
**VIENES?**

# Challenges

- Sizes, fonts, orientations
- Languages
- Scale
- Efficiency

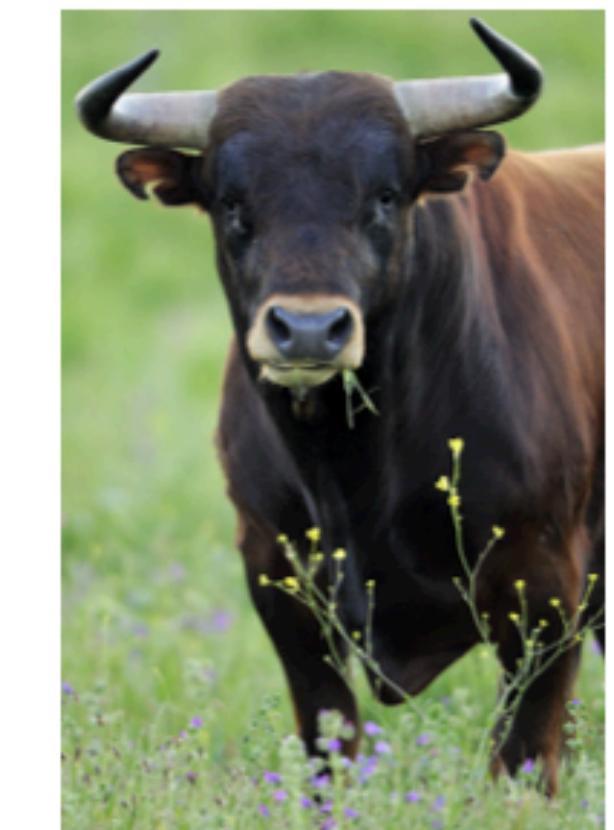
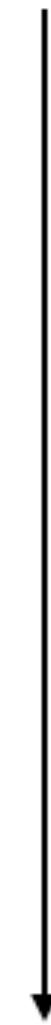


Classes for image

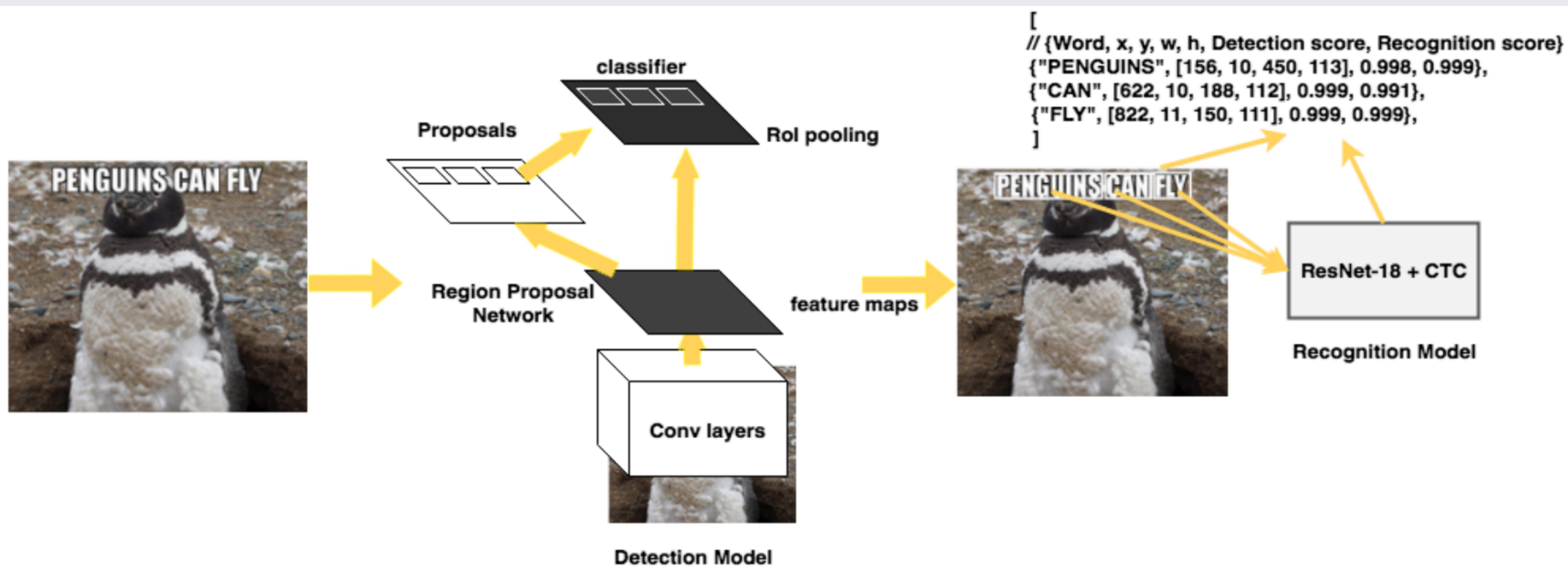


bulldog    dog    petshop    pet    english  
puppy    petfood    bull    dogfood    creche  
microchip    haldol    doggie    peludo    kennel

**bulldog**

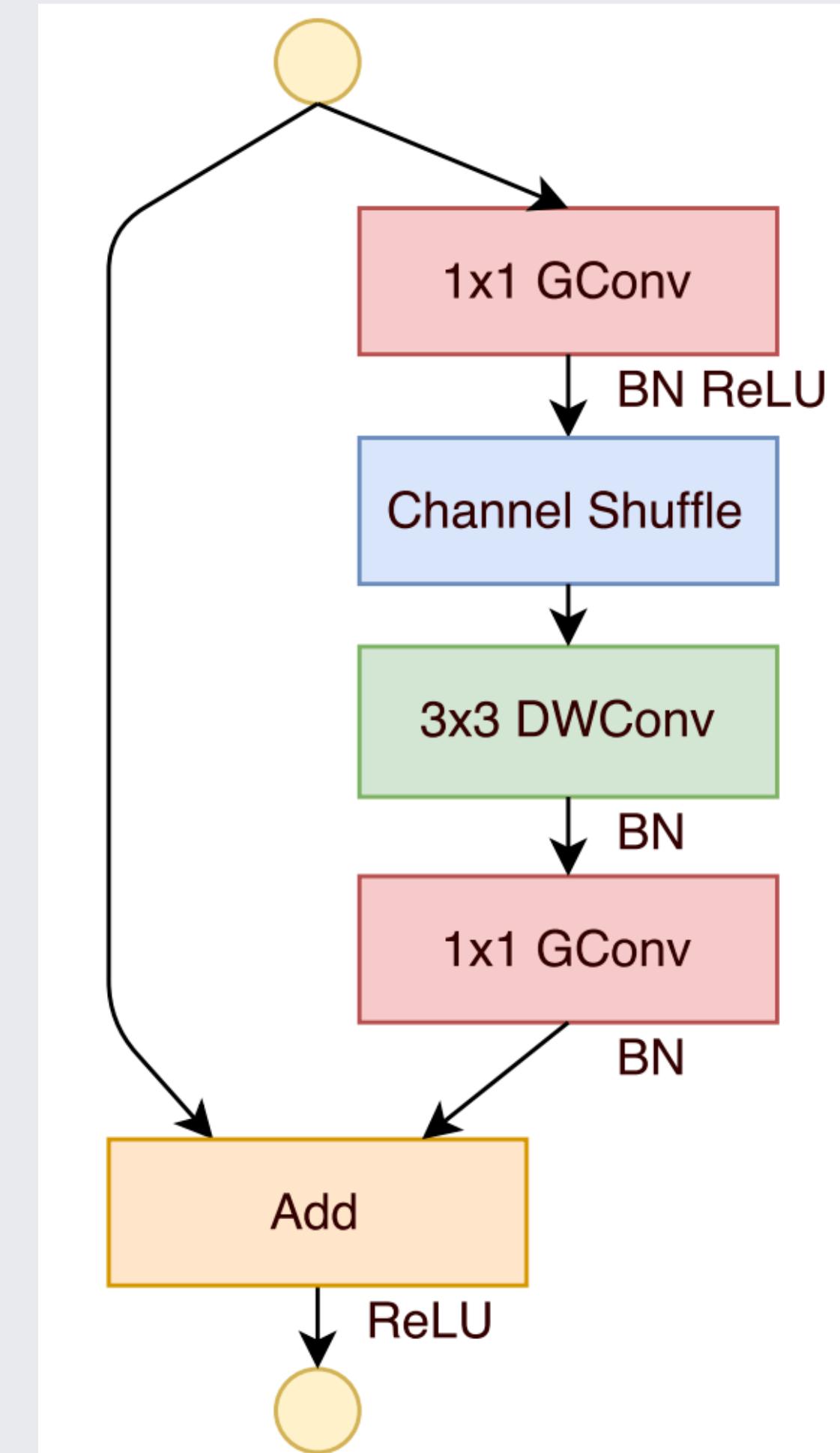
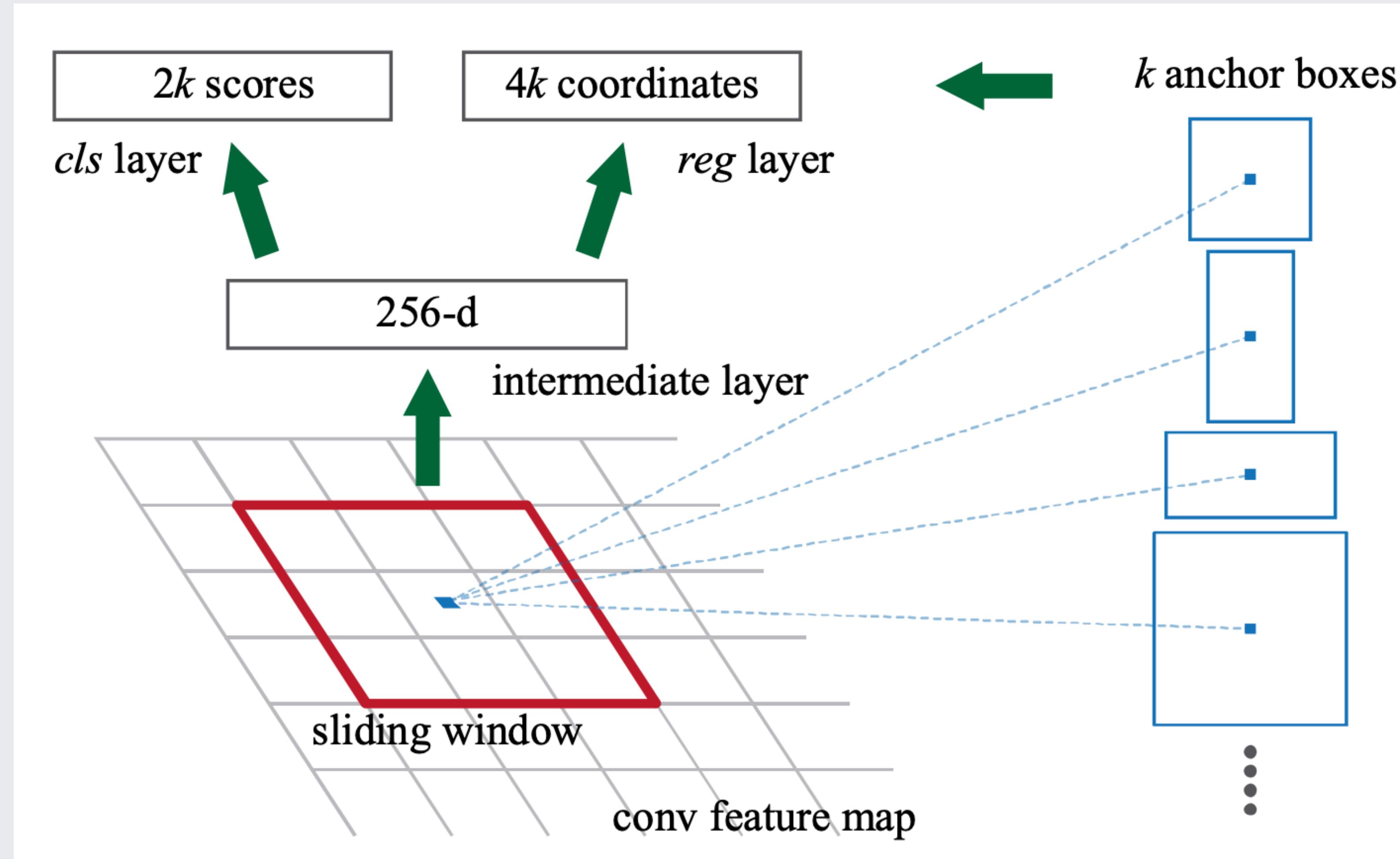


# Architecture



Large scale system for text detection and recognition in images, KDD 2018,  
Viswanath Sivakumar, Albert Gordo, Fedor Borisyuk

# Text Detection



Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks, Ren et al.

ShuffleNet: An Extremely Efficient Convolutional Neural Network for Mobile Devices, Zhang et al.

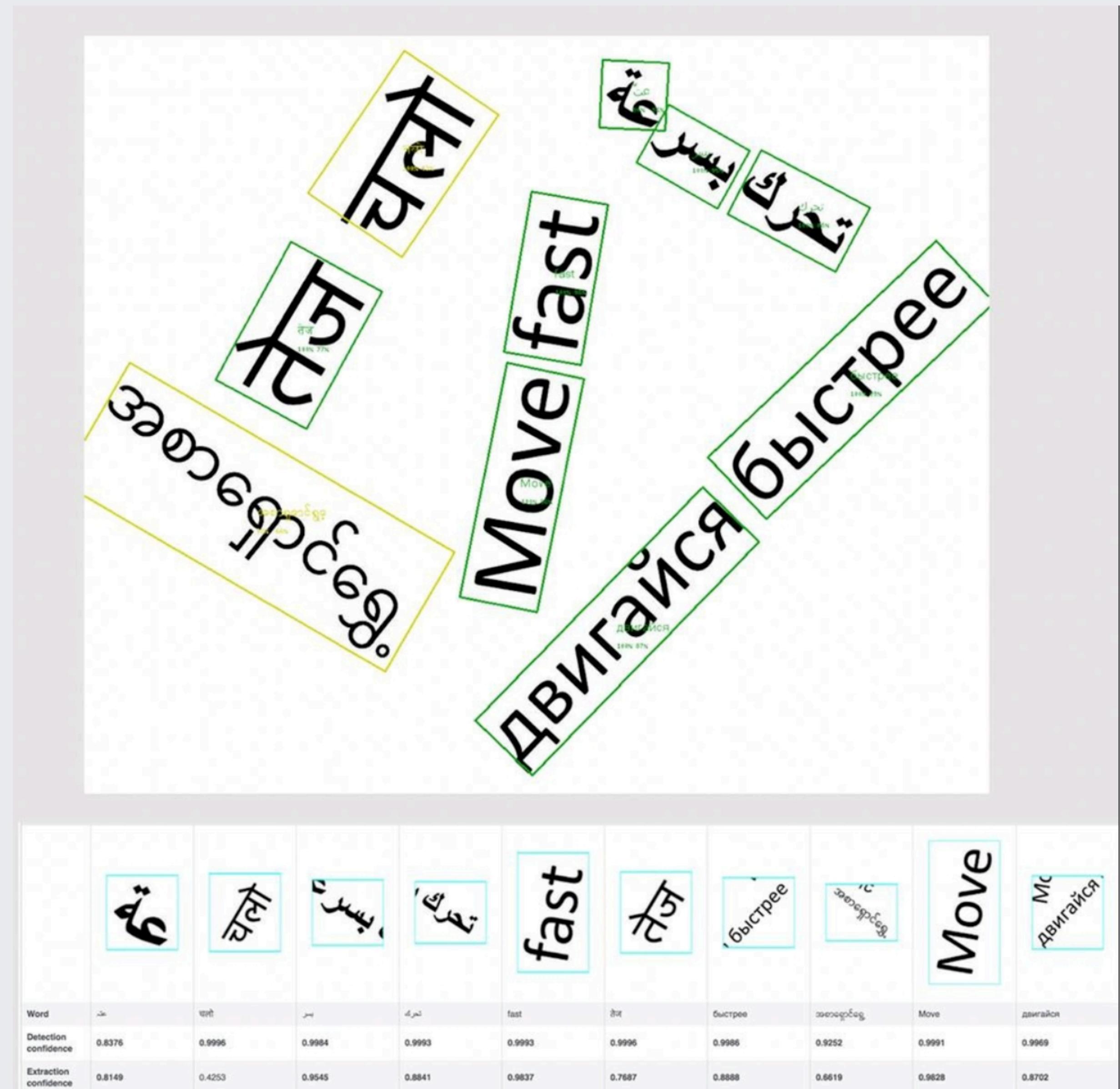
# Orientations

The diagram illustrates the detection and orientation analysis of text in an image. The text 'Move fast' and its Russian translation 'двигайся быстрее' are shown in various orientations. Colored boxes (red, green, yellow) highlight specific words or parts of words, indicating their detected orientation and confidence.

Below the diagram is a table showing word-level detection and extraction confidence scores:

Word	کی	لے	لے	لے	لے	fast	سے	Move	بے	میں	دھیاں
Detection confidence	0.7507	0.8731	0.9578	0.8830	0.9578	0.9006	0.9402	0.9900	0.9842		
Extraction confidence	0.7739	0.7489	0.6879	0.8583	0.7650	0.2286	0.4658	0.4614	0.2924		

# Orientations



# Improving Rotated Text Detection with Rotation Region Proposal Networks, Huang et al.

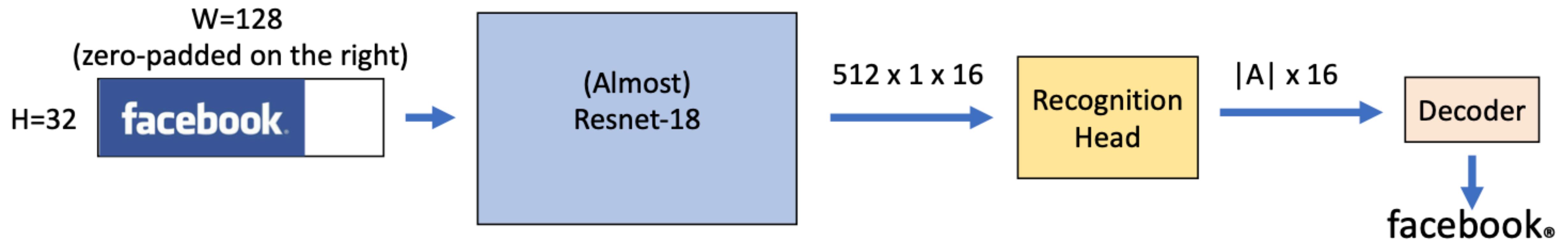
# Text Recognition



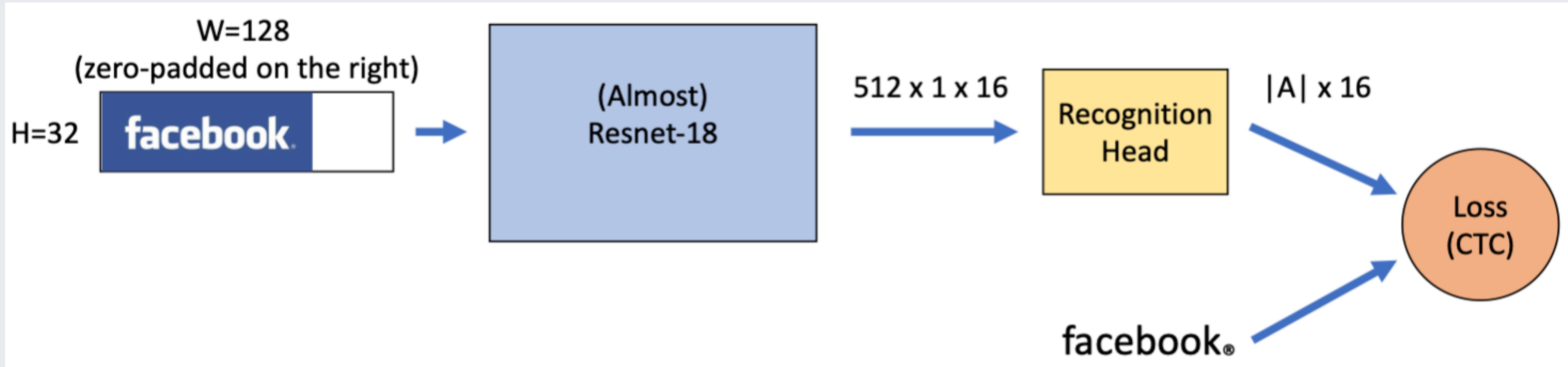
# Text Recognition

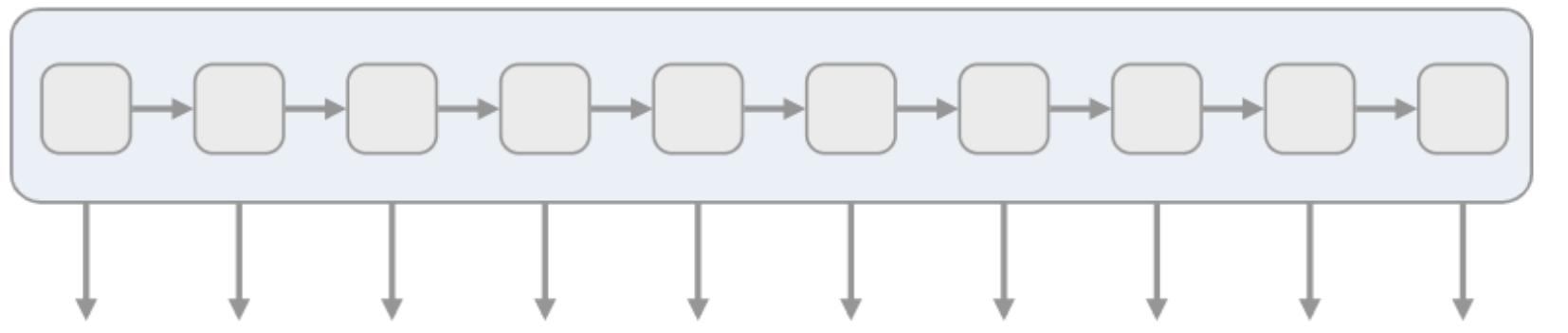


# Text Recognition



# Text Recognition





h	h	h	h	h	h	h	h	h	h	h
e	e	e	e	e	e	e	e	e	e	e
o	o	o	o	o	o	o	o	o	o	o
€	€	€	€	€	€	€	€	€	€	€

The input is fed into an RNN,  
for example.

The network gives  $p_t(a | X)$ ,  
a distribution over the outputs  
 $\{h, e, |, o, \epsilon\}$  for each input step.

h	e	€			€			o	o
h	h	e			€	€		€	o
€	e	€			€	€		o	o

With the per time-step output  
distribution, we compute the  
probability of different sequences

h	e			o
e			o	
h	e		o	

By marginalizing over alignments,  
we get a distribution over outputs.

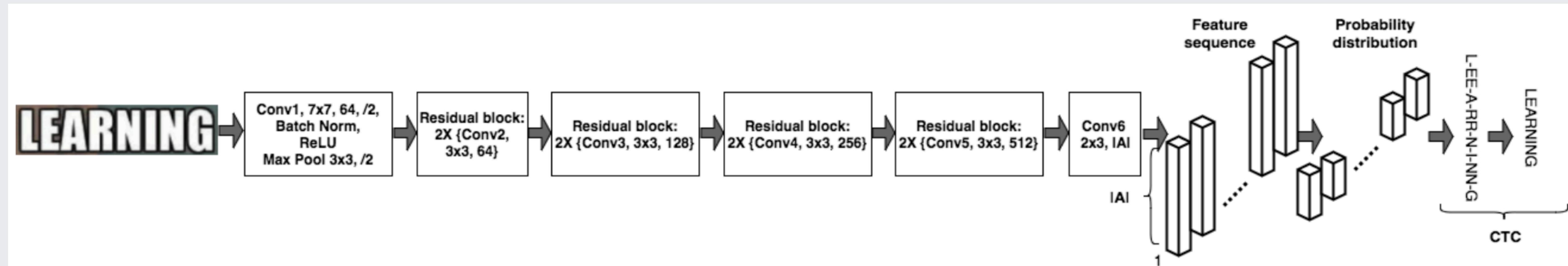
$$p(Y | X) = \sum_{A \in \mathcal{A}_{X,Y}} \prod_{t=1}^T p_t(a_t | X)$$

The CTC  
conditional  
**probability**

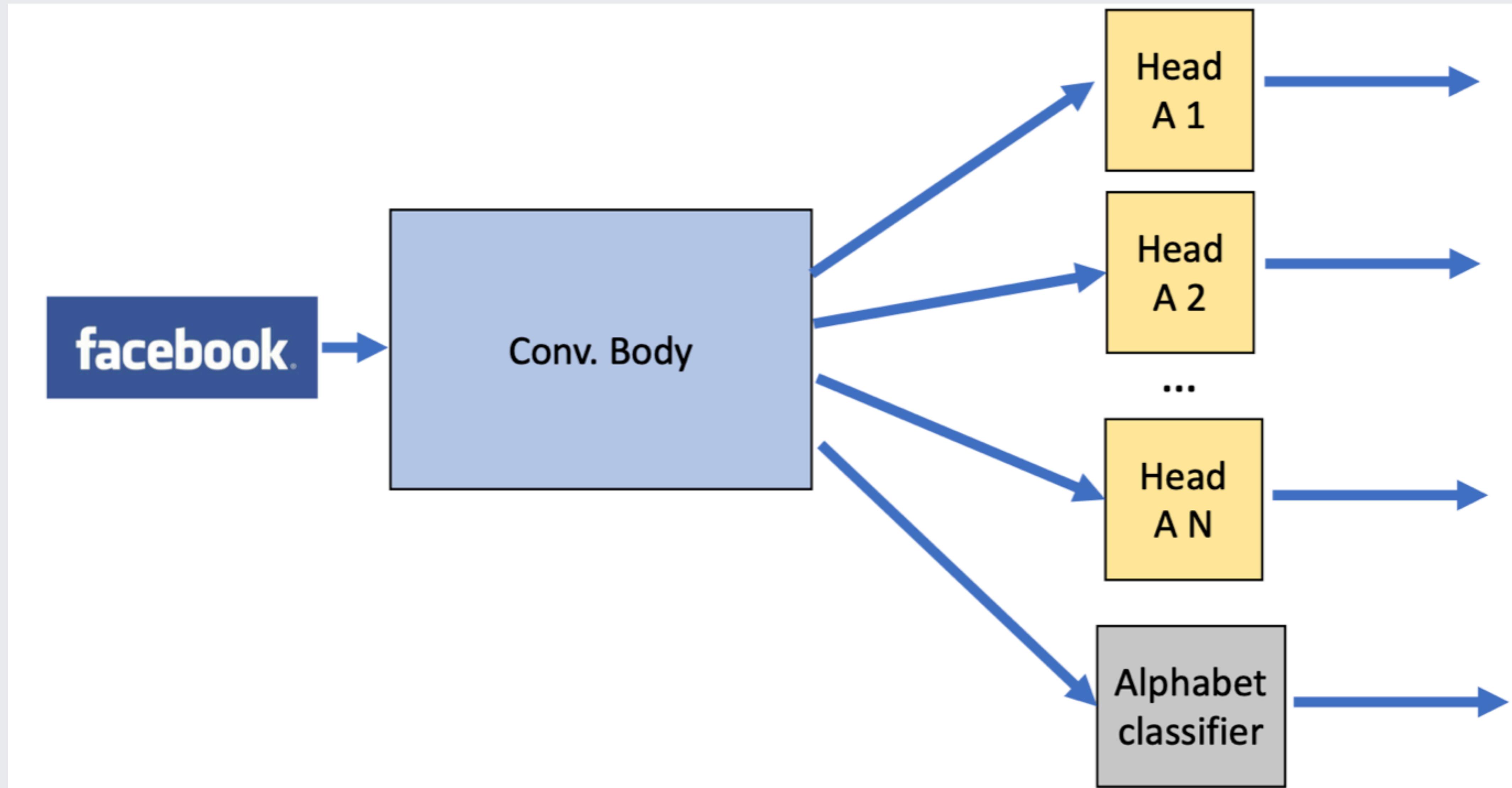
**marginalizes**  
over the set of  
valid  
alignments

computing the  
**probability** for a single  
alignment step-by-  
step.

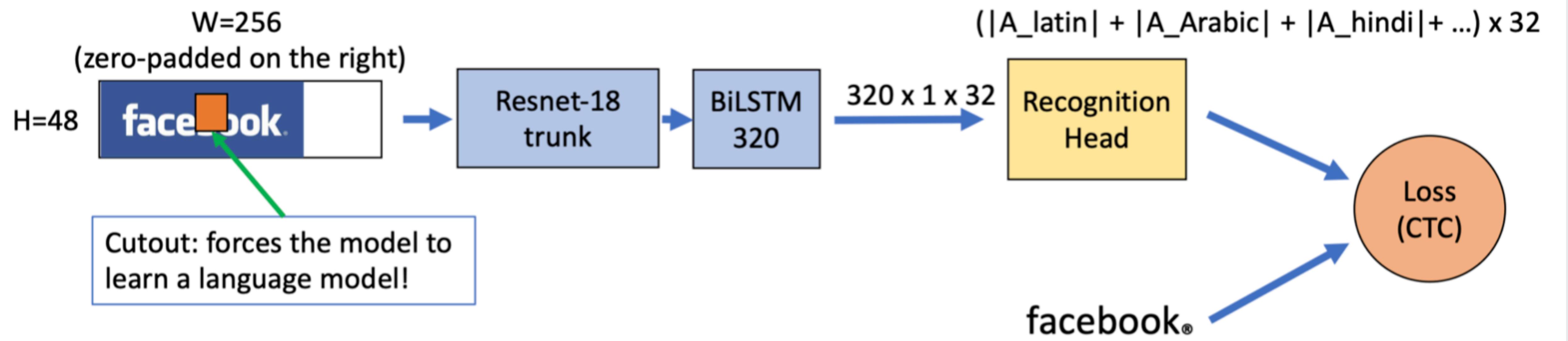
# Text Recognition



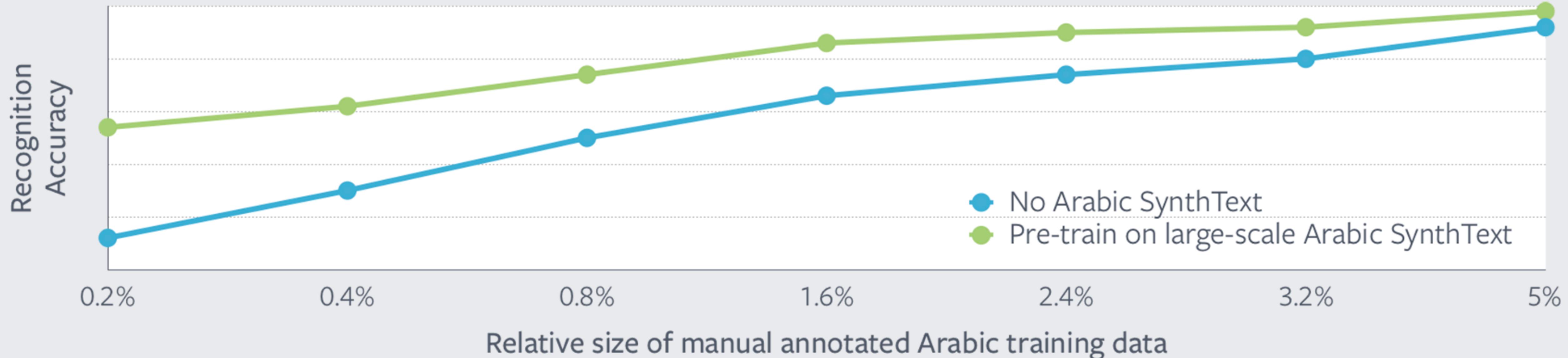
# Multilingual



# Multilingual



# Synthetic Data



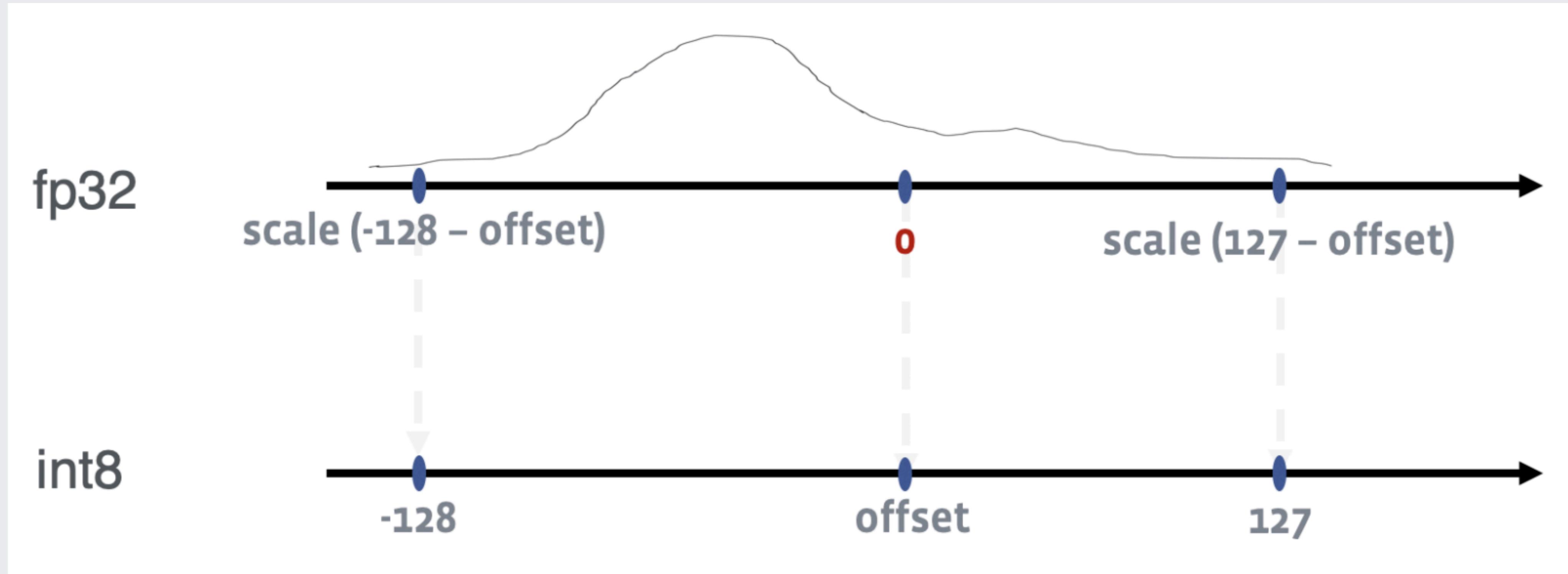
# Inference

1B images/day x 5 sec/image = Lots of servers!

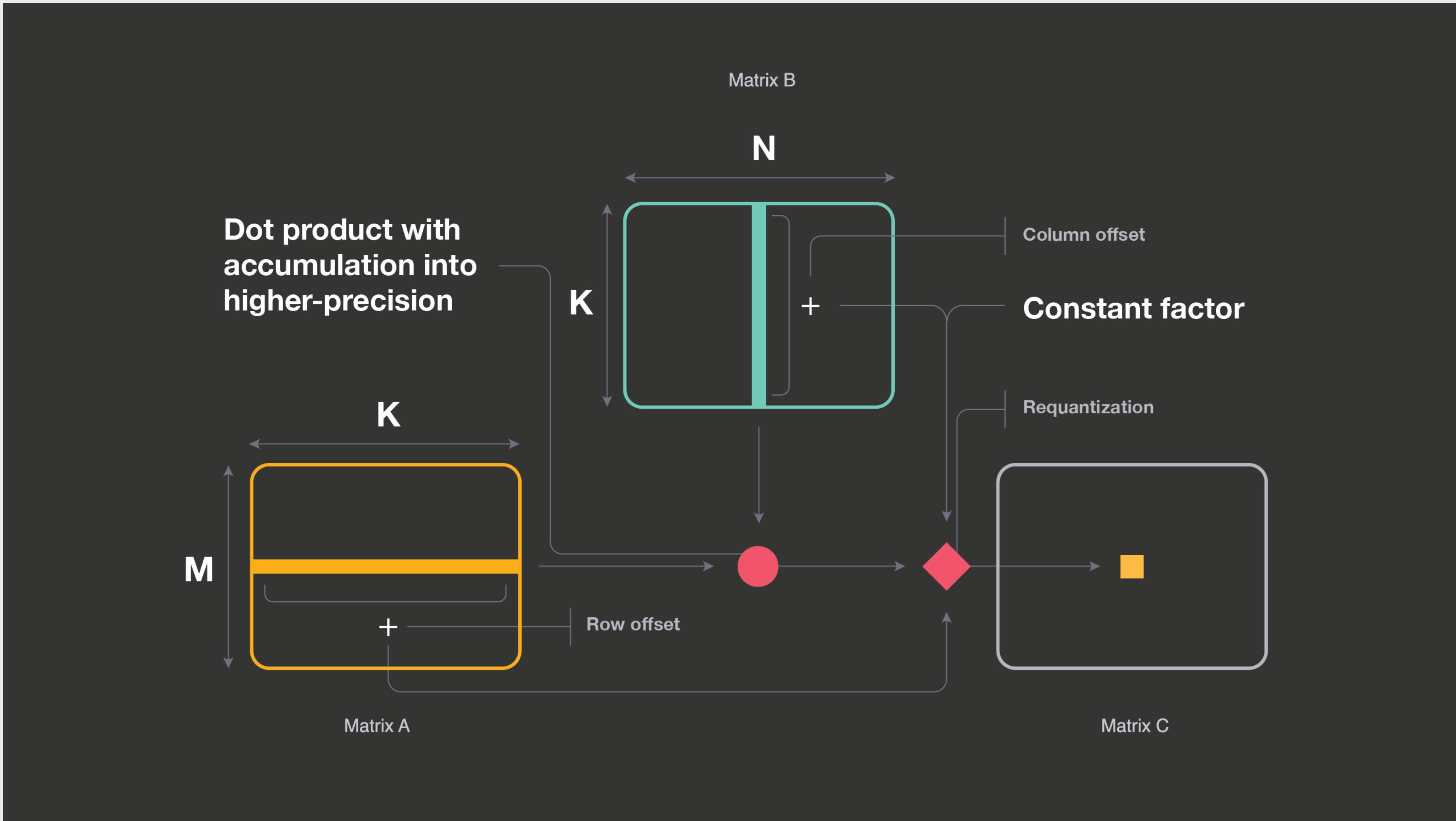
# Int8 Quantization

**Quantization:**  $x_q = \text{clip}(\text{round}(x/\text{scale}) + \text{offset}, -128, 127)$

**De-quantization:**  $x = \text{scale} \times (x_q - \text{offset})$



# Int8 Quantization



# But what about accuracy?

- No quantization error for 0
- Fuse Convolution and ReLU
- L2 Error Minimization vs Min-Max
- Don't quantize the first layer

Initial accuracy gap: 5%

After: 0.2%

# Int8 Quantization

Great for memory bandwidth bound systems!

But what about CPU?

Intel Architecture	fp32 GEMM with 32-bit accumulation	int8 GEMM with 32-bit accumulation	int8 GEMM with 16-bit accumulation
Broadwell	32 ops/cycle	32 ops/cycle	64 ops/cycle
Skylake	64 ops/cycle	64 ops/cycle	192 ops/cycle

16-bit accumulation is faster!  
But what about accuracy drop due to overflows?

# Outlier-Aware Quantization

A block of unsigned  
int8 activation  
matrix X

0	10	0	10
30	200	10	250
40	0	0	0
0	20	10	30

A block of signed  
int8 weight matrix  
 $W^T$

-30	20	30	10
-10	<b>-120</b>	<b>120</b>	10
0	20	30	-20
-20	<b>-120</b>	<b>120</b>	10

x

=

A partial sum of  
block of output Y in  
signed int16

-300	-2400	2400	200
-7900	<b>-53200</b>	<b>55200</b>	4600
-1200	800	1200	400
-800	-5800	6300	300

# Outlier-Aware Quantization

A block of unsigned int8 activation matrix X

0	10	0	10
30	200	10	250
40	0	0	0
0	20	10	30

Non-outlier of  $W^T$

-30	20	30	10
-10	<b>0</b>	<b>0</b>	10
0	20	30	-20
-20	<b>0</b>	<b>0</b>	10

A partial sum of block of output Y in signed int16

-300	0	0	200
-7900	<b>800</b>	<b>1200</b>	4600
-1200	800	1200	400
-800	200	300	300

+ int32 accumulation

A block of unsigned int8 activation matrix X

0	10	0	10
30	200	10	250
40	0	0	0
0	20	10	30

Outlier of  $W^T$  (typically very sparse)

0	0	0	0
0	<b>-120</b>	<b>120</b>	0
0	0	0	0
0	<b>-120</b>	<b>120</b>	0

A partial sum of block of output Y in signed int32

0	-2400	2400	0
0	<b>-54000</b>	<b>54000</b>	0
0	0	0	0
0	-6000	6000	0

$$W = \text{Dense}_W + \text{Sparse}_W$$

$$W: [-128, 127]$$

$$\text{Dense}_W: [-64, 63]$$

$$\text{Sparse}_W: \text{Outliers}$$

$$X * W =$$

$$X * \text{Dense}_W \text{ (16-bit acc)}$$

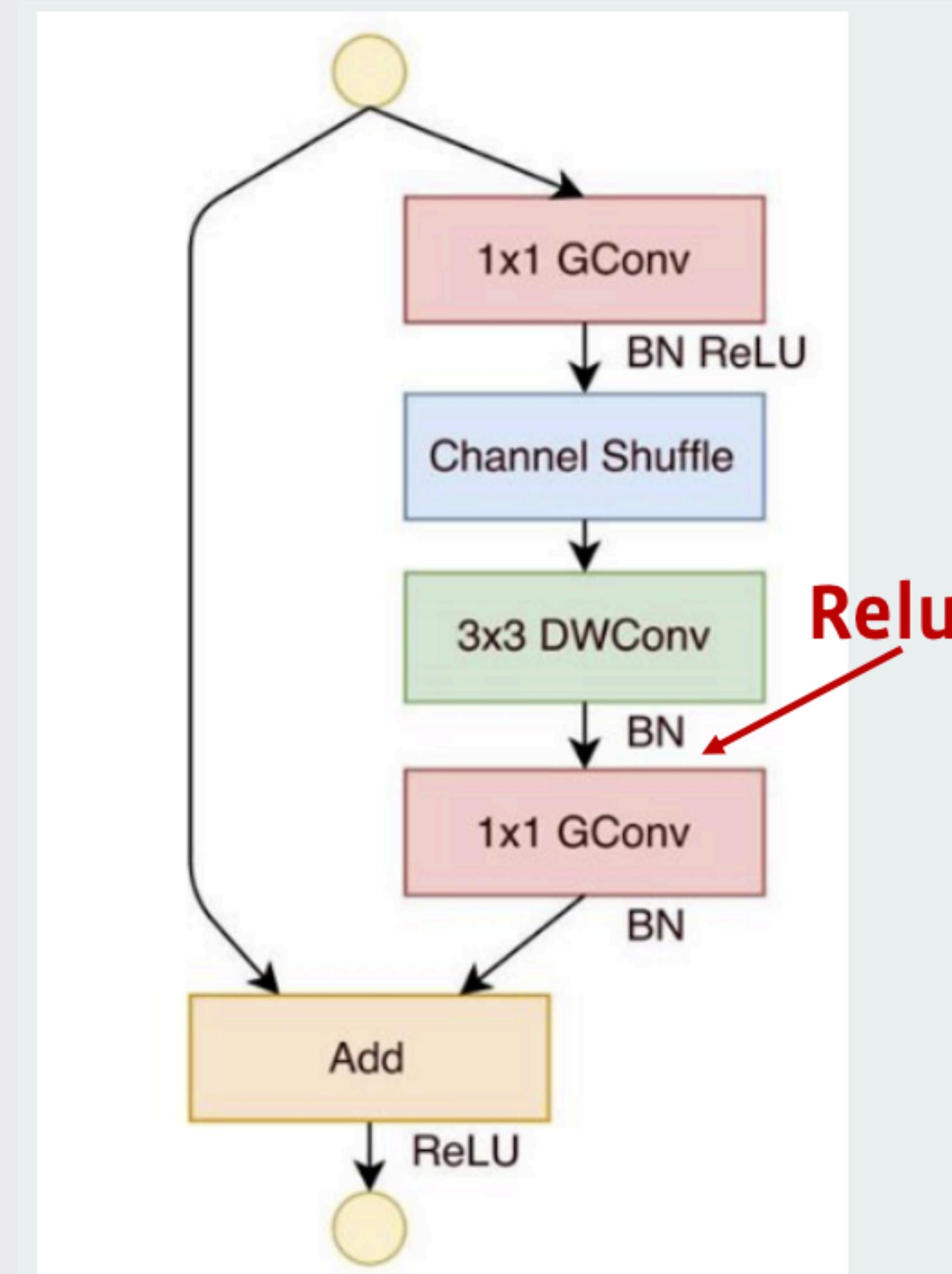
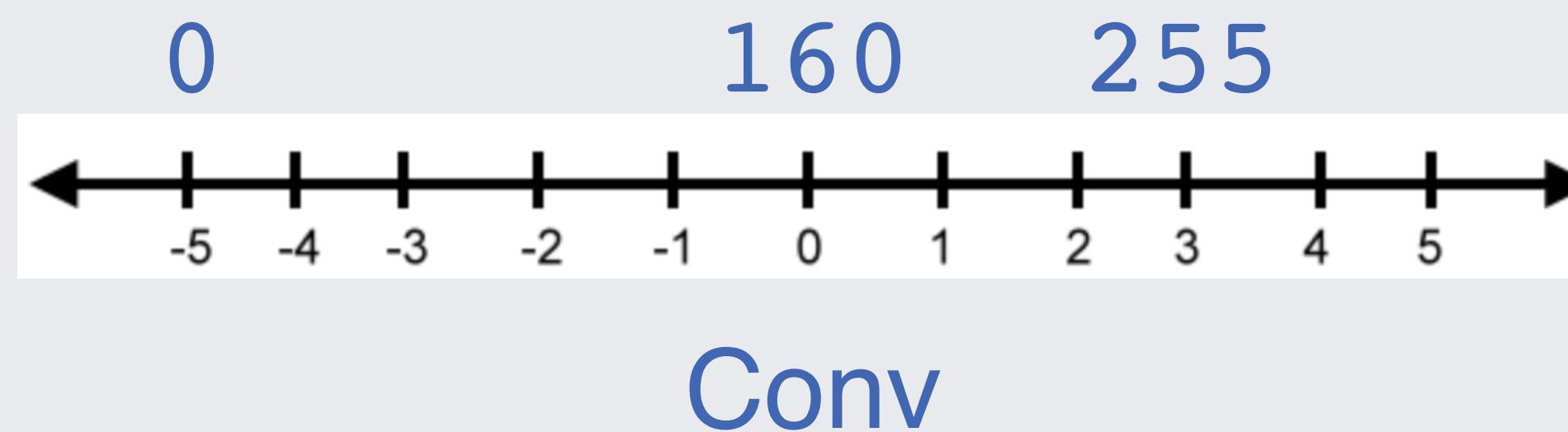
+

$$X * \text{Sparse}_W \text{ (32-bit acc)}$$

# Model Co-Design

Outlier-aware quantization

- Int8 Quantization with 16-bit accumulation
- Further CPU speedup

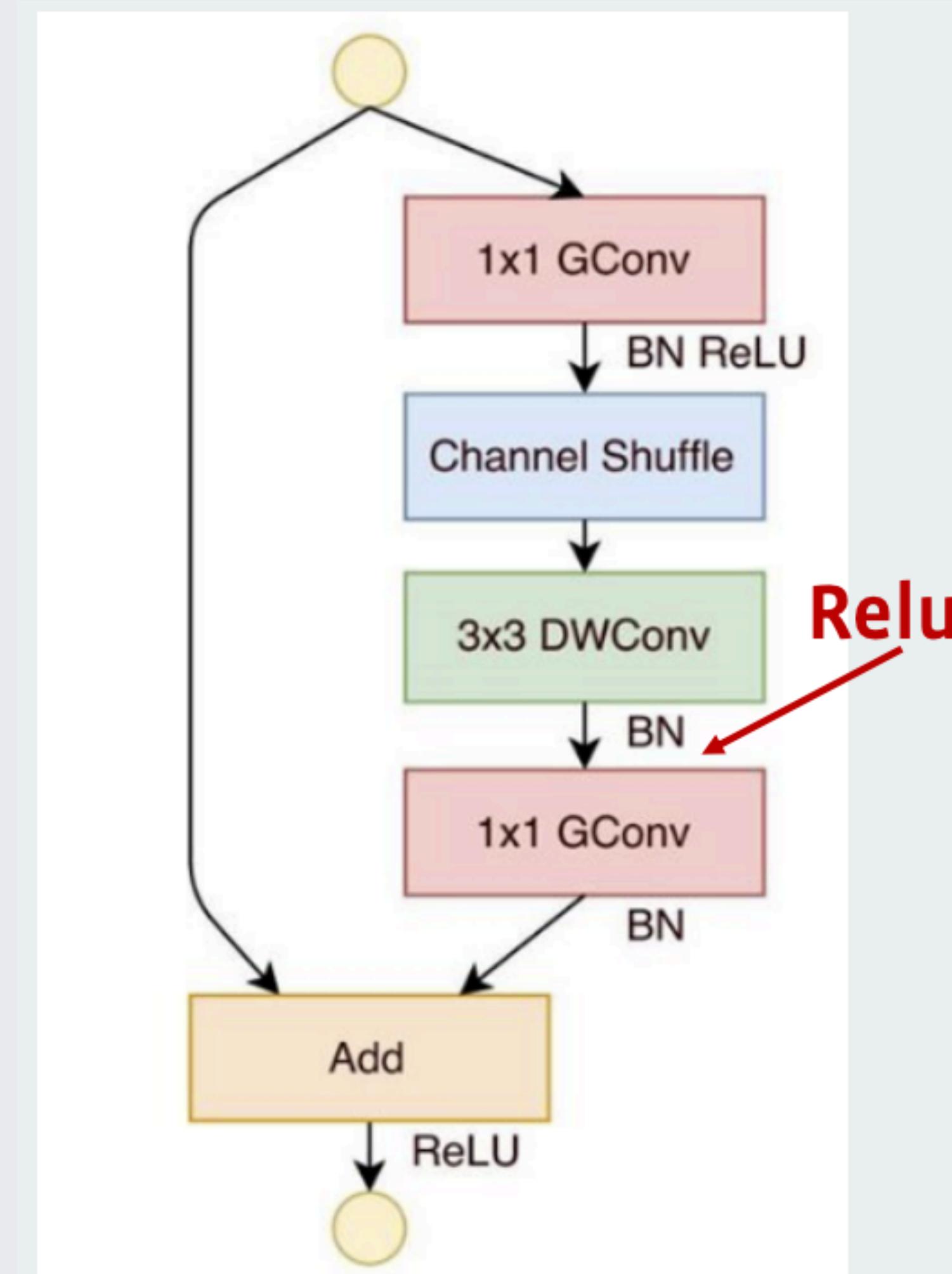
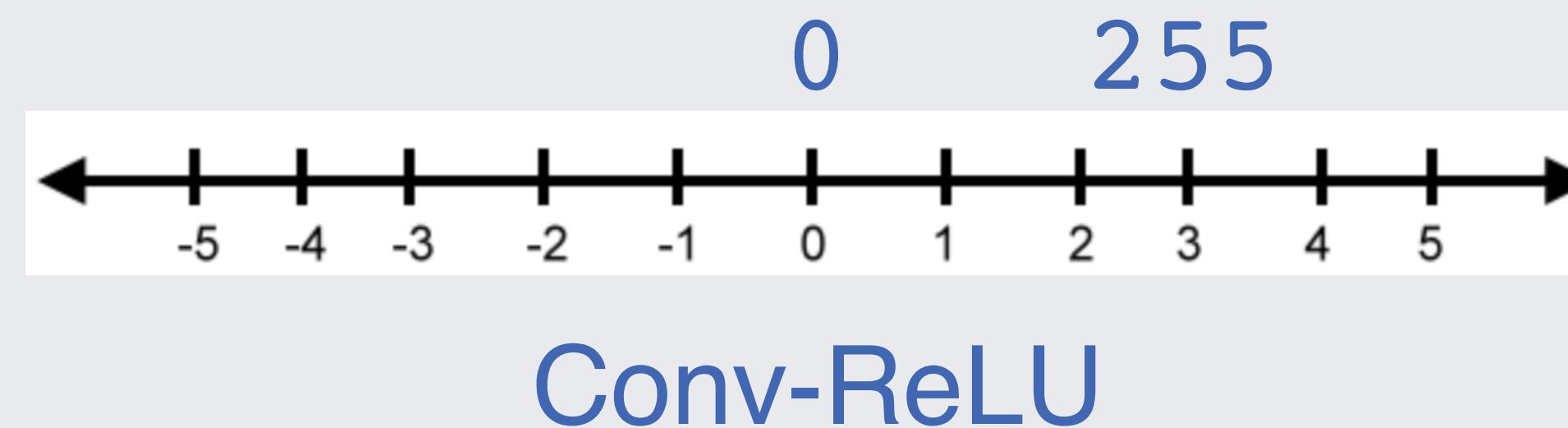


ShuffleNet

# Model Co-Design

Outlier-aware quantization

- Int8 Quantization with 16-bit accumulation
- Further CPU speedup

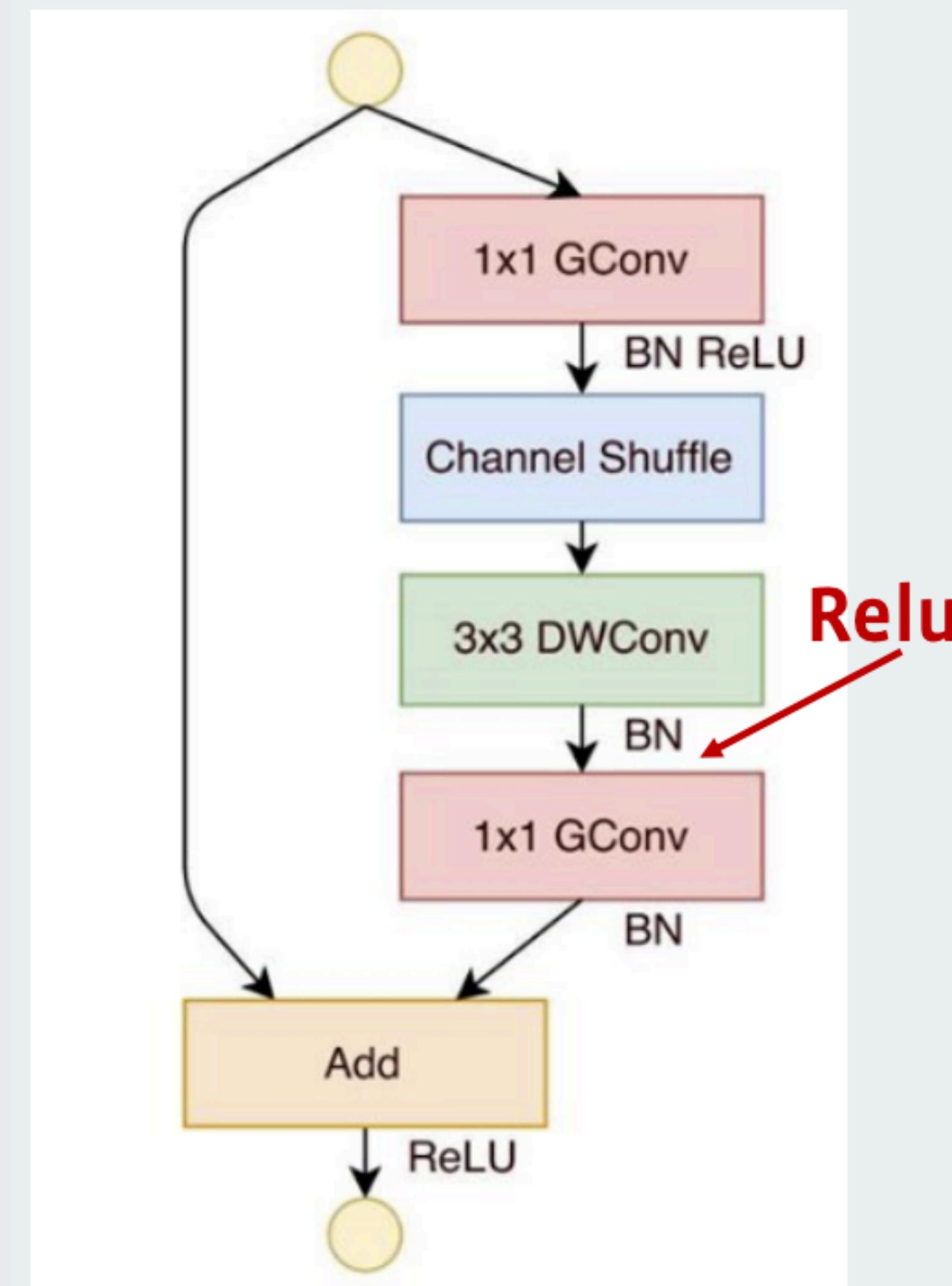
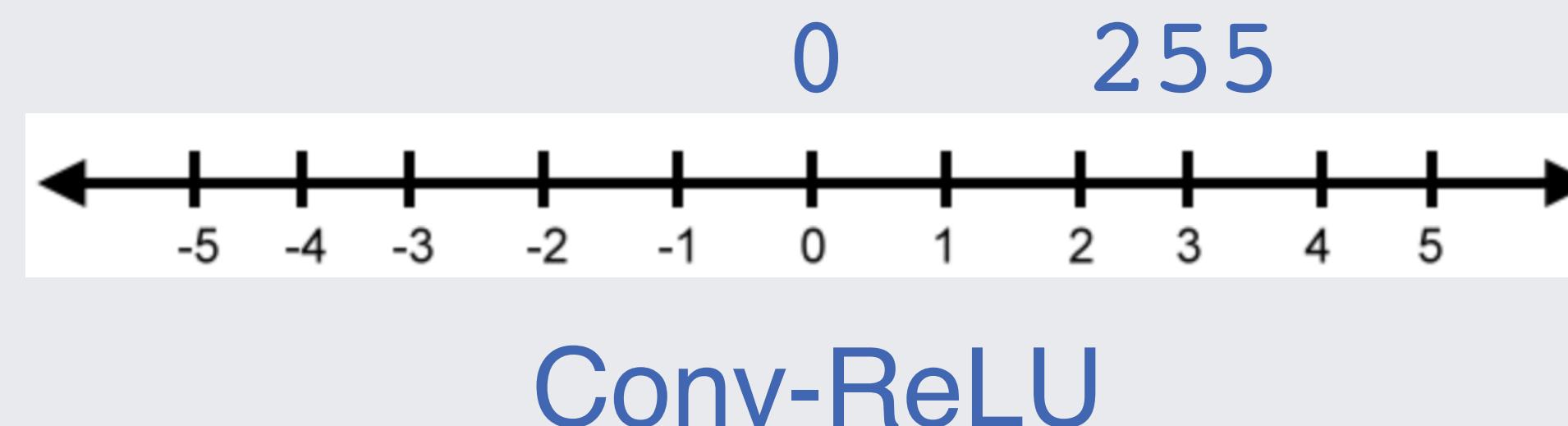


ShuffleNet

# Model Co-Design

Outlier-aware quantization

- Int8 Quantization with 16-bit accumulation
- Further CPU speedup



Co-design your model and efficiency optimizations together!

ShuffleNet

# Int8 Quantization



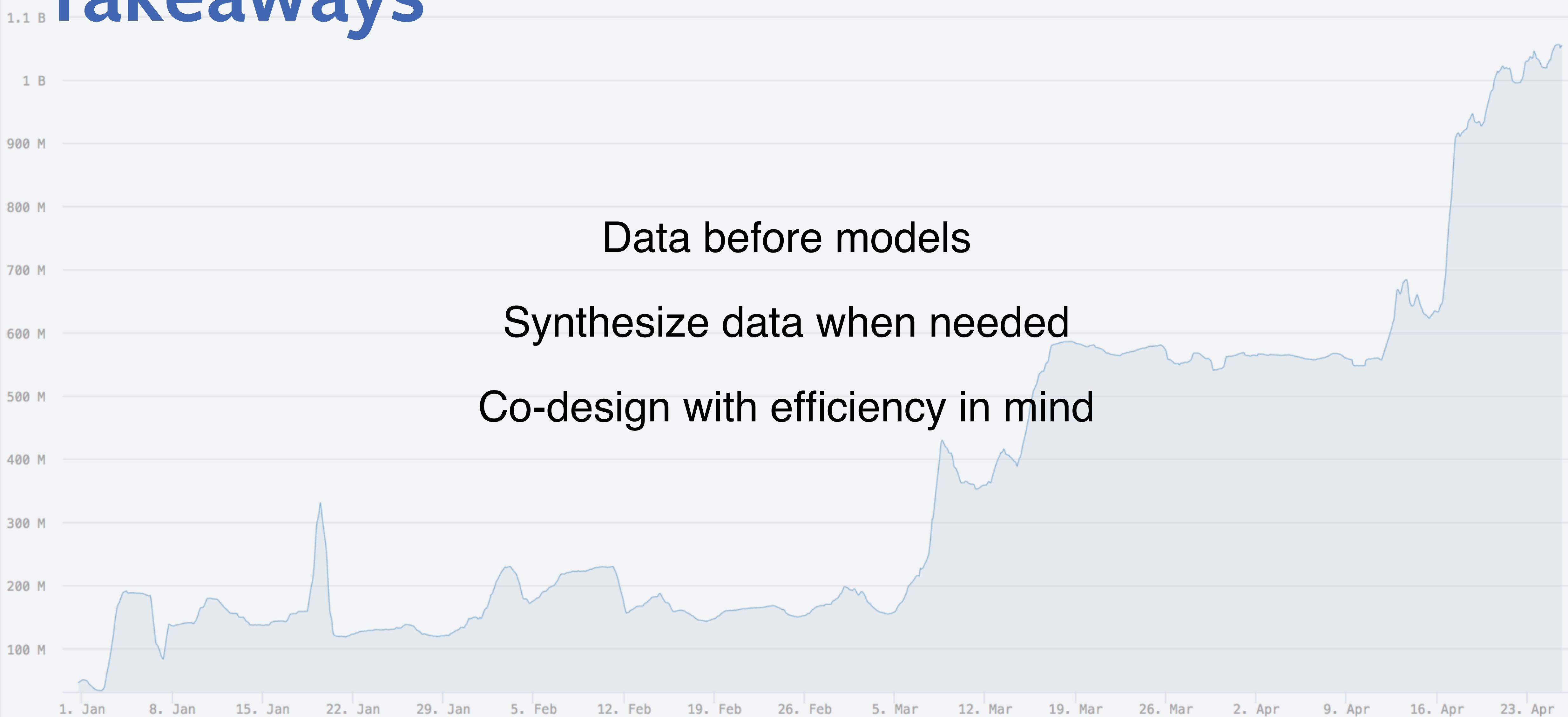
2.4x images/server

2x faster

[github.com/pytorch/FBGEMM](https://github.com/pytorch/FBGEMM)



# Takeaways



# Takeaways

Data before models

Synthesize data when needed

Co-design with efficiency in mind

Thank You!

