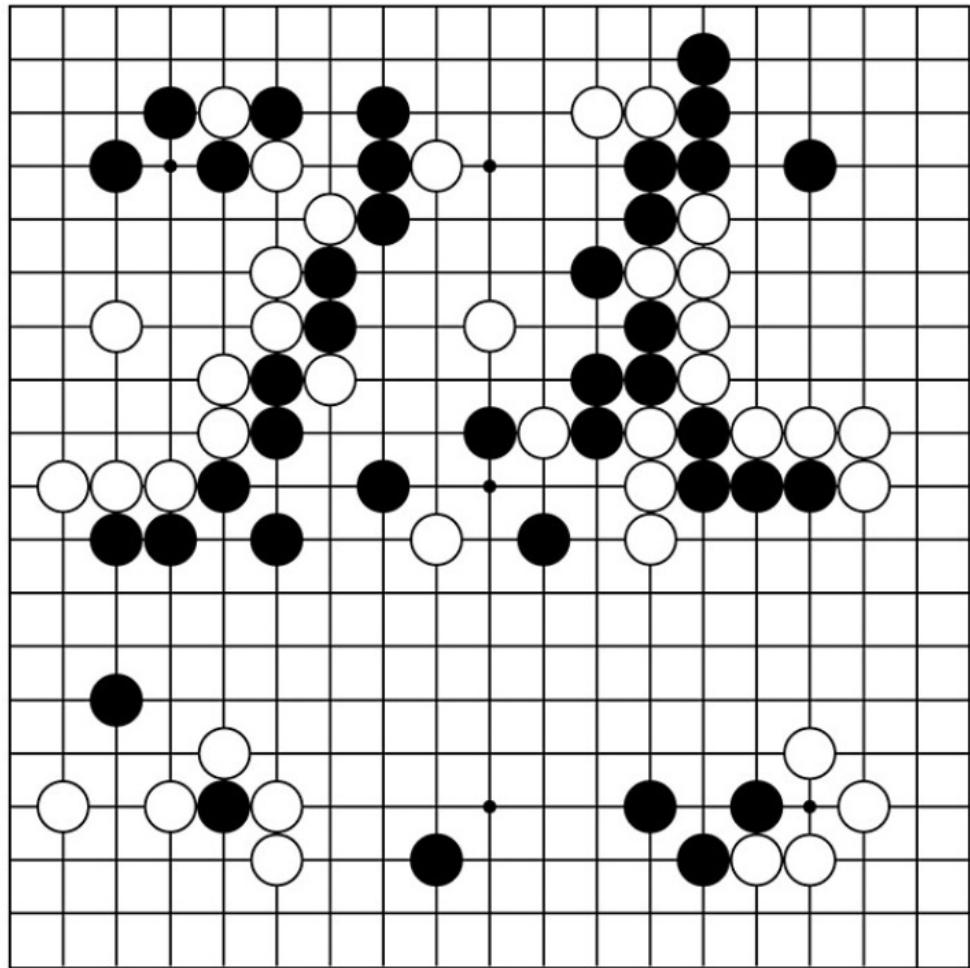


Scalable Visual Pre-training: Past, Present, and Paths Forward

Saining Xie
Courant CS, NYU
saining.xie@nyu.edu



3^{361} possible states? 🤔

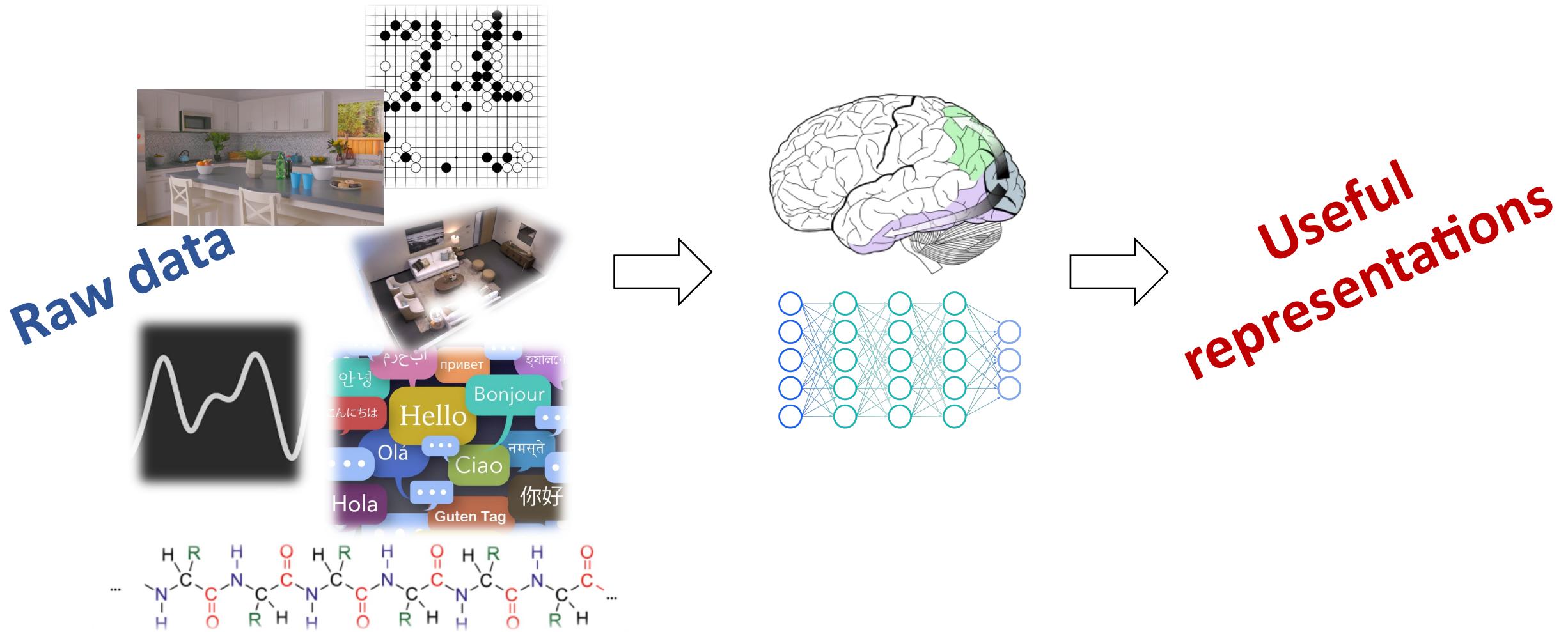
(>> number of atoms in the universe)

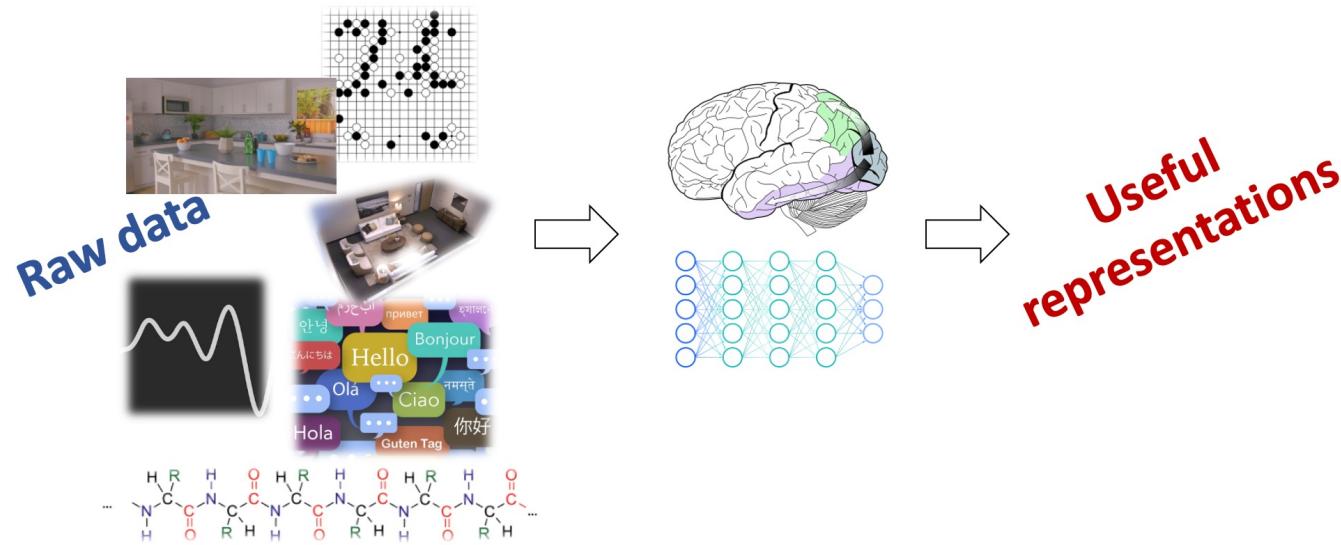
Figure: AlphaGo vs Lee Sedol G4 move 78



$256^3 \times 600 \times 800$
possible states?

Figure: detectron2





- In some cases, we can directly “*use*” the useful representations.
- In many applications, we need to “***transfer***” the useful representations to some downstream tasks.

Transfer learning

- tasks: same
- labels: source task

Transductive

different domains

Domain
Adaptation

different languages

Cross-lingual
learning

- tasks: different
- labels: target task

Inductive

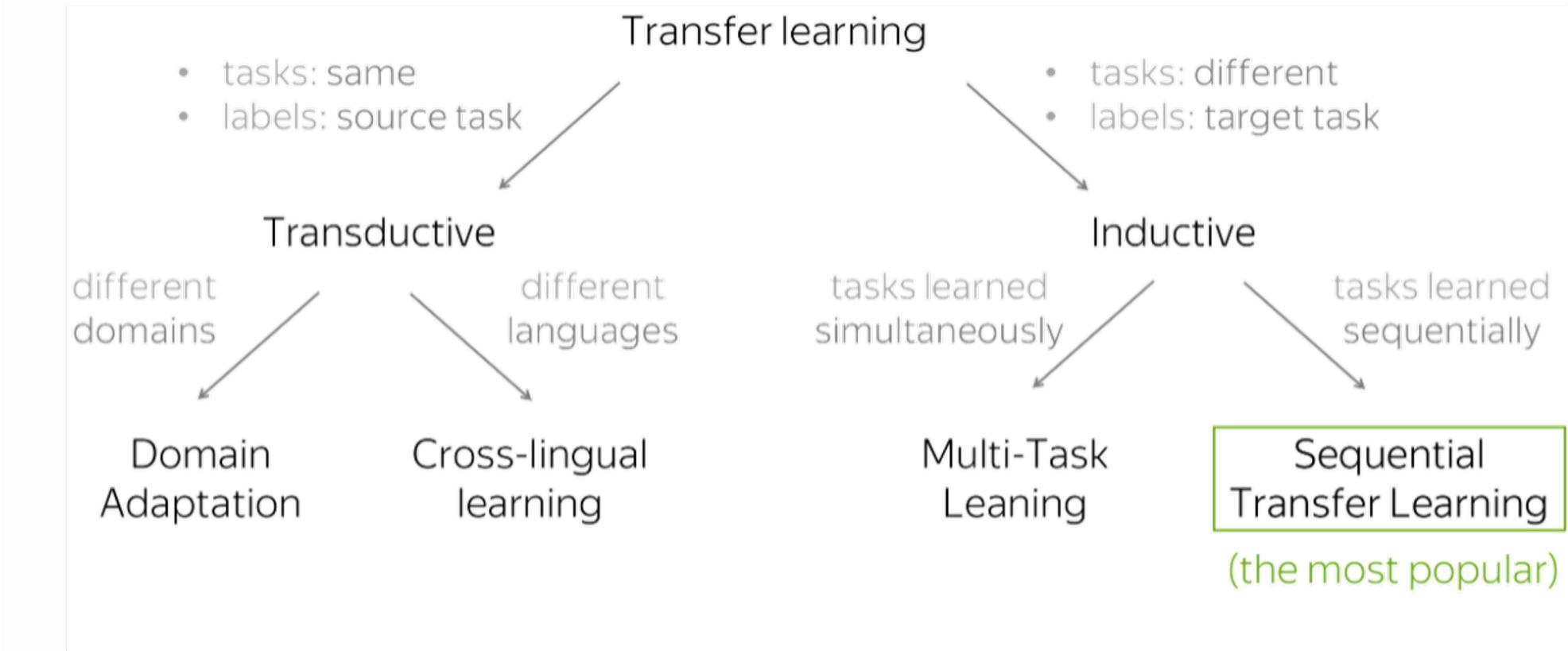
tasks learned
simultaneously

Multi-Task
Learning

tasks learned
sequentially

Sequential
Transfer Learning
(the most popular)

This taxonomy is from [Sebastian Ruder's blog post](#).



Pre-training → Fine-tuning*

Architecture: what to train?

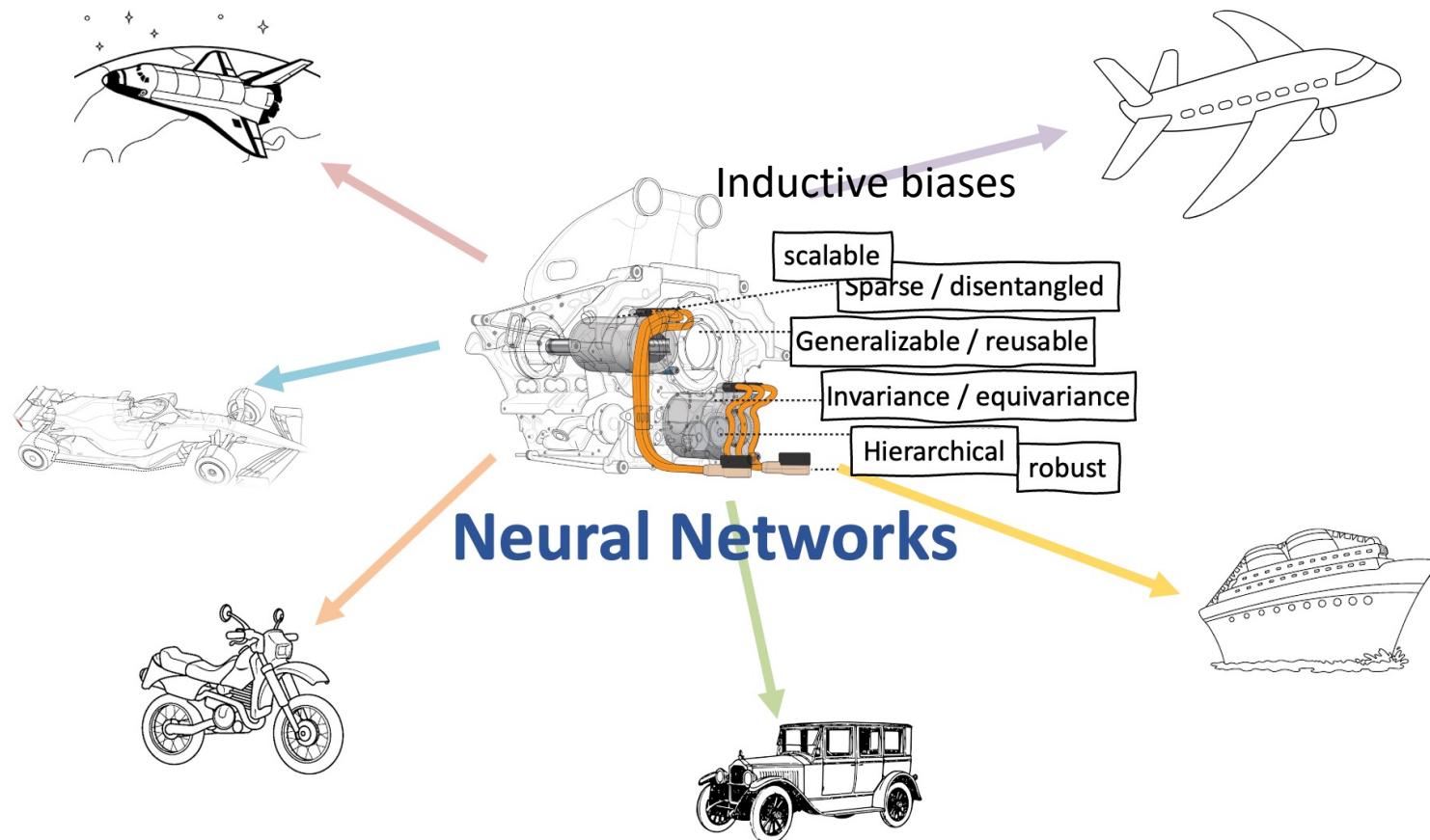
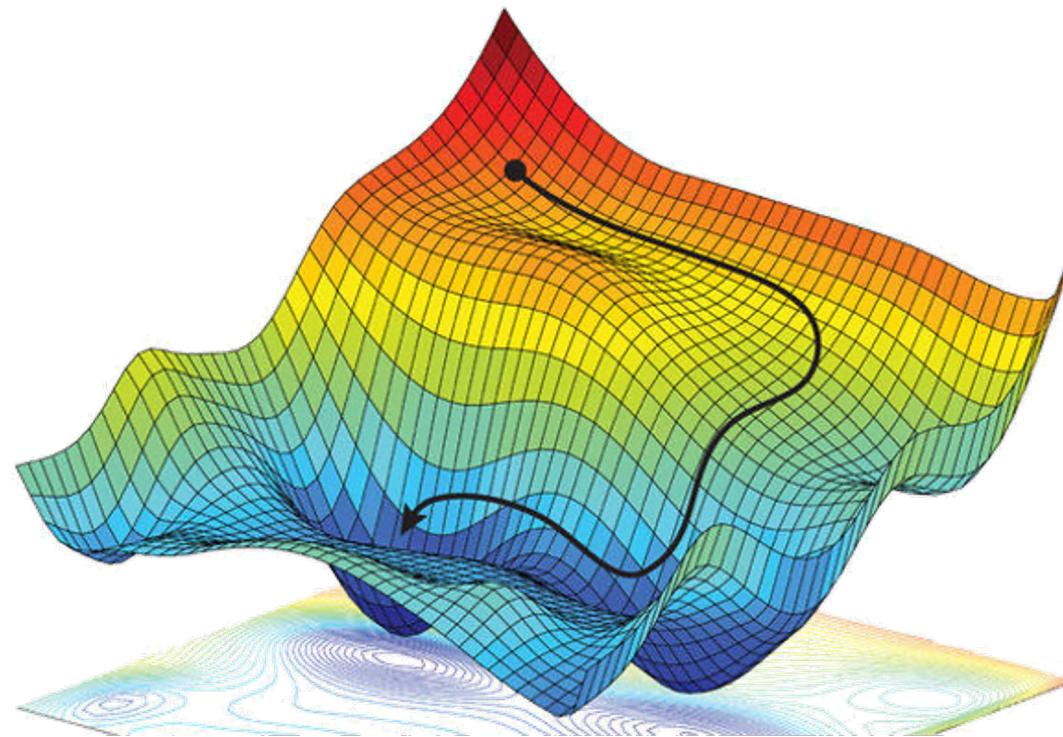
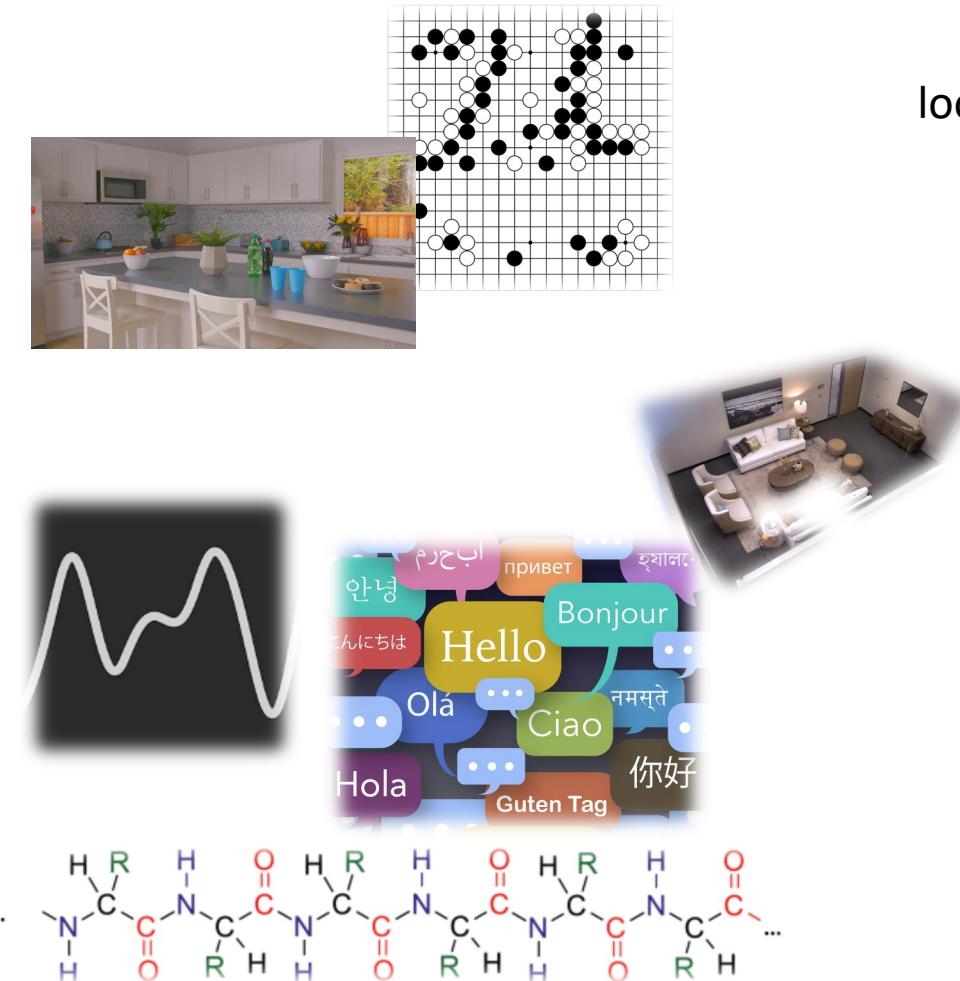


Illustration: Craig Scarborough, clipart-library.com

Objective: how to train?



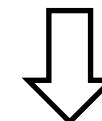
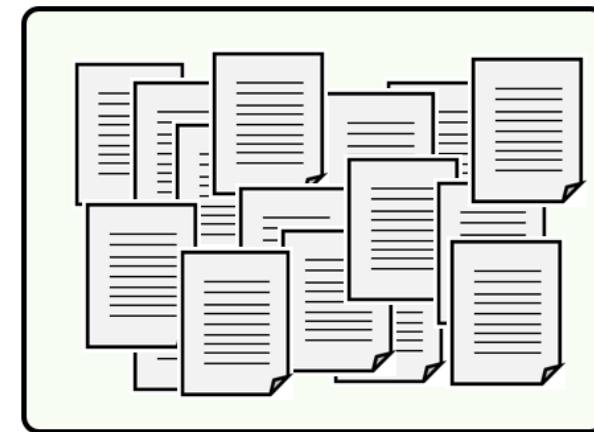
Data: where to train?



loosely organized

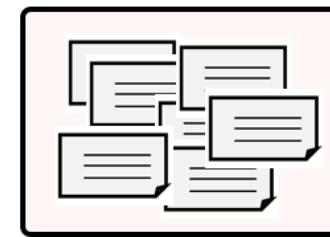
huge diverse corpus

general domain



not huge, or not diverse, or both

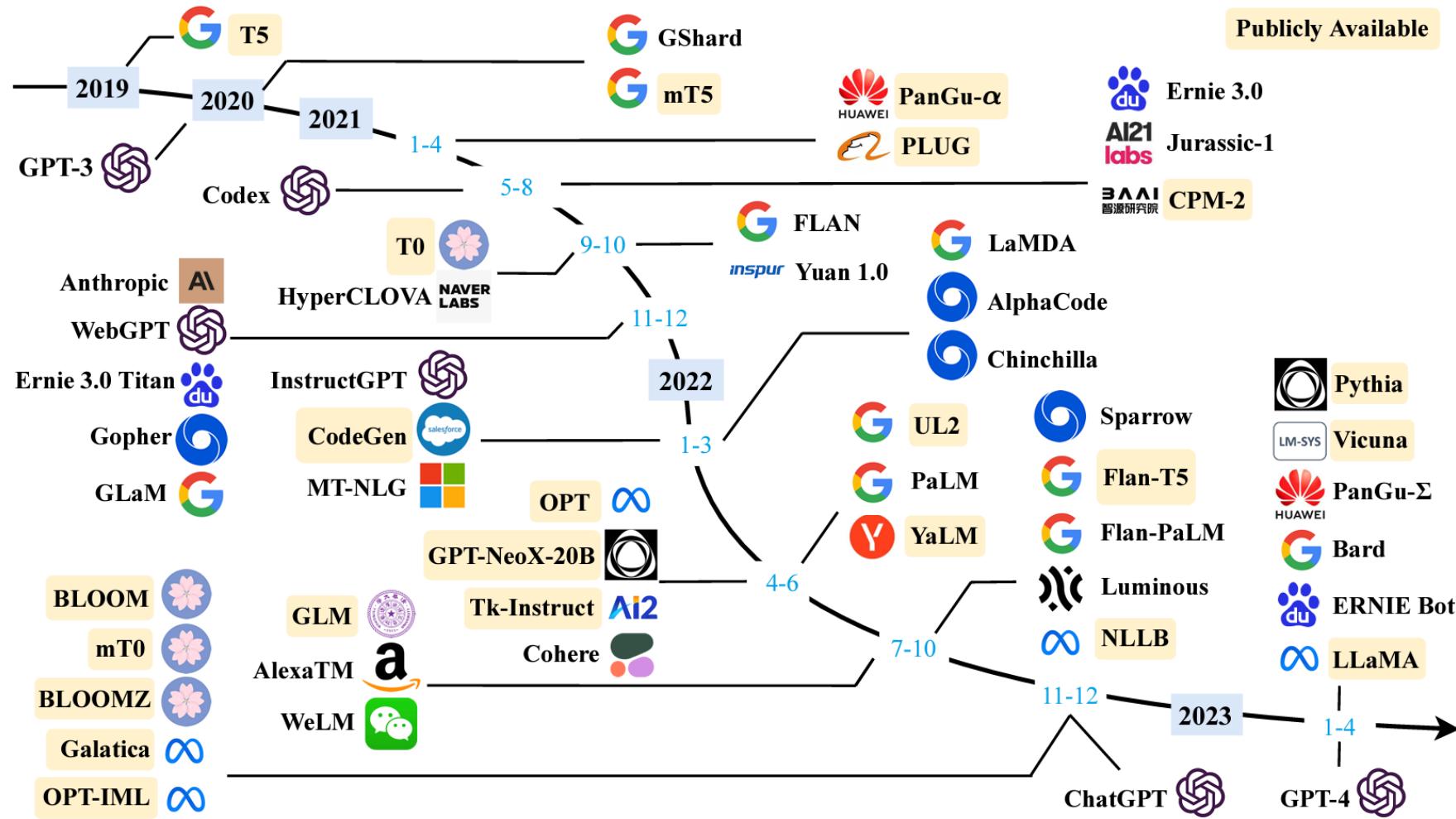
carefully curated



task-specific
domain

task-specific
annotations

Large pre-trained language models

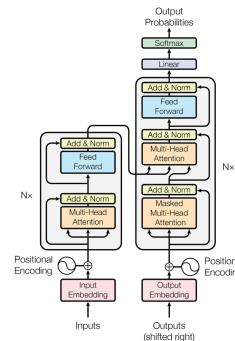


Large pre-trained language models

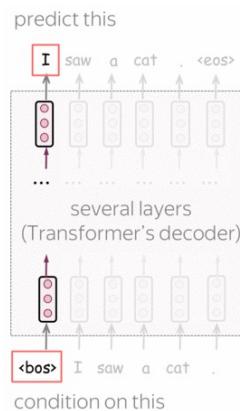
Data:



Architecture:



Objective:



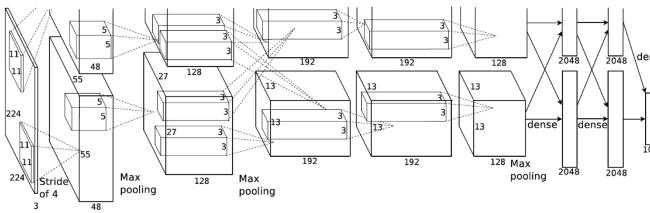
Visual Pre-training has a longer history...

Data:

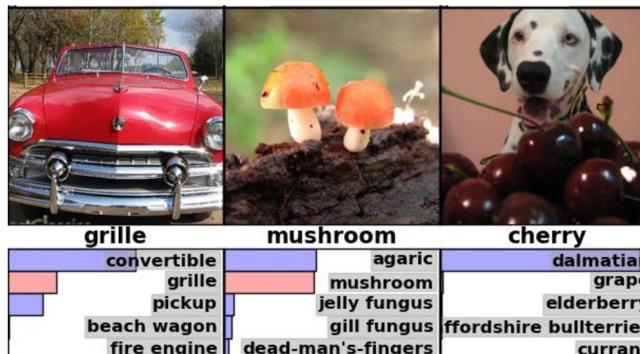


[Deng et al., 2009]

Architecture:



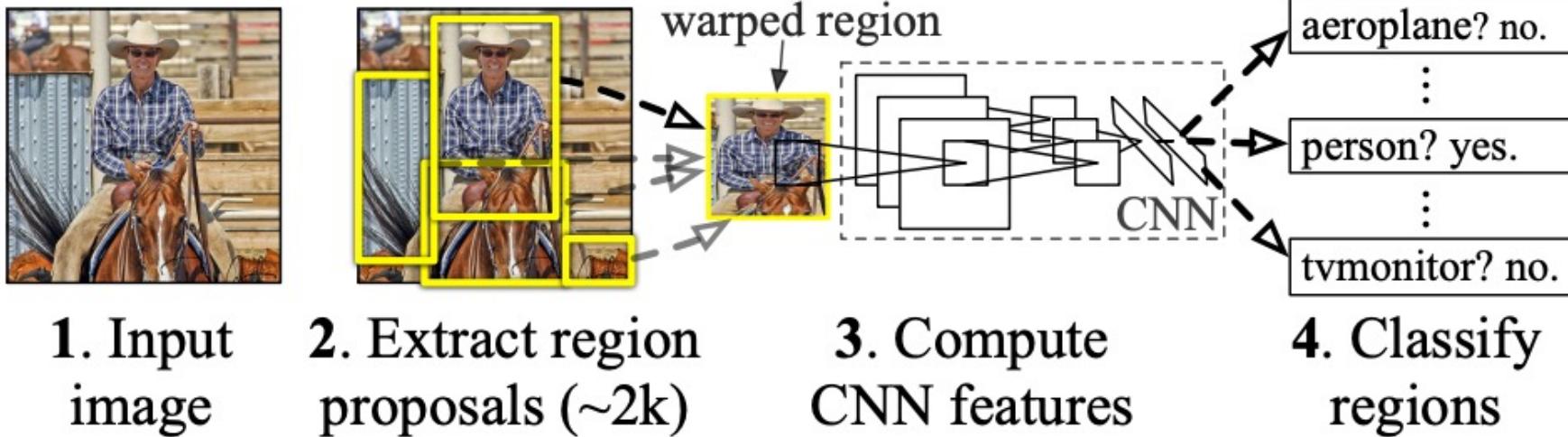
Objective:



Visual Pre-training has a longer history...

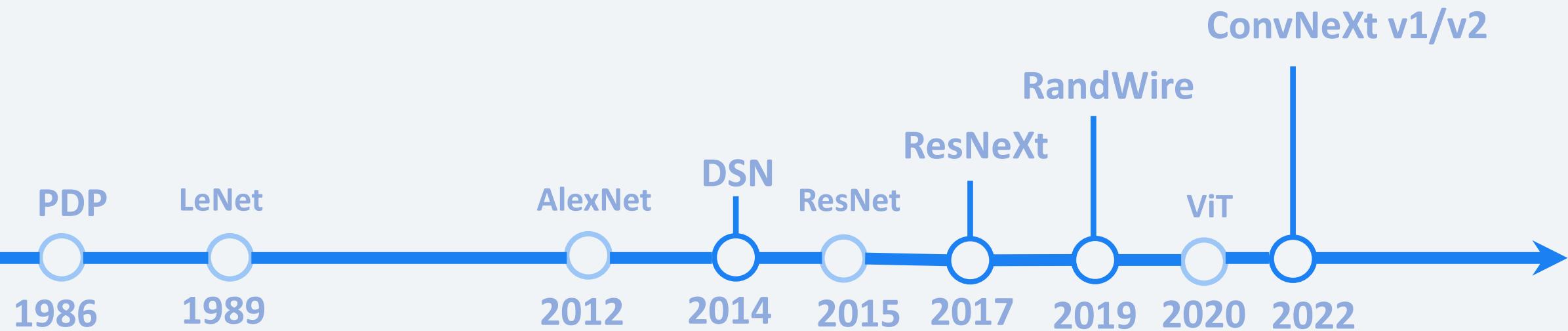
Arguably the first major success of “pre-training”...

R-CNN: *Regions with CNN features*



Architecture / Objective / Data

1. How to design neural network architectures



Architecture / Objective / Data

2. Training objectives beyond supervised classification: Are labels necessary?



Architecture / Objective / **Data**

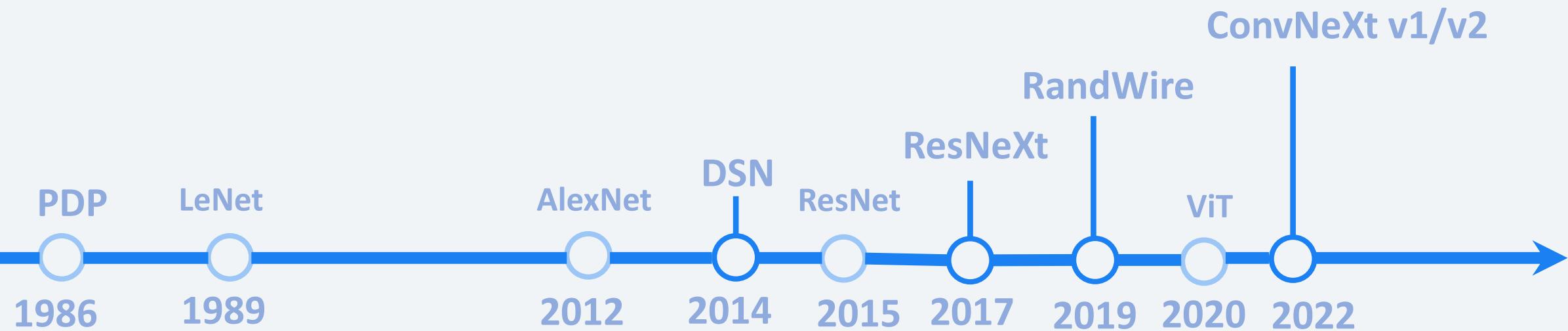
3. What data to use for visual pre-training

**Multi-Modality
Noisy Labels**



Architecture / Objective / Data

1. How to design neural network architectures



Connectionism

[A general framework for parallel distributed processing. Rumelhart et al., 1986]

PARALLEL DISTRIBUTED PROCESSING

Explorations in the Microstructure of Cognition

Volume 1: Foundations

David E. Rumelhart James L. McClelland
and the PDP Research Group

Chisato Asanuma
Francis H. C. Crick
Jeffrey L. Elman
Geoffrey E. Hinton
Michael I. Jordan

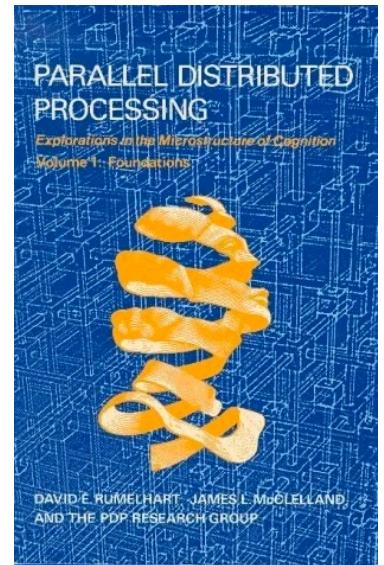
Alan H. Kawamoto
Paul W. Munro
Donald A. Norman
Daniel E. Rabin
Terrence J. Sejnowski

Paul Smolensky
Gregory O. Stone
Ronald J. Williams
David Zipser

Institute for Cognitive Science
University of California, San Diego

(PDP group is now at Stanford)

- A *set of processing units*
- A *state of activation*
- An *output function* for each unit
- A *pattern of connectivity* among units
- A *propagation rule* for propagating patterns of activities through the network of connectivities
- An *activation rule* for combining the inputs impinging on a unit with the current state of that unit to produce a new level of activation for the unit.
- A *learning rule* whereby patterns of connectivity are modified by experience
- An *environment* within which the system must operate



1986

Convolutional Neural Networks

[Learning Internal Representations by Error Propagation. Rumelhart et al., 1986]

ConvNet using BP

- Receptive field
- Translation equivariance
- Trained by error propagation

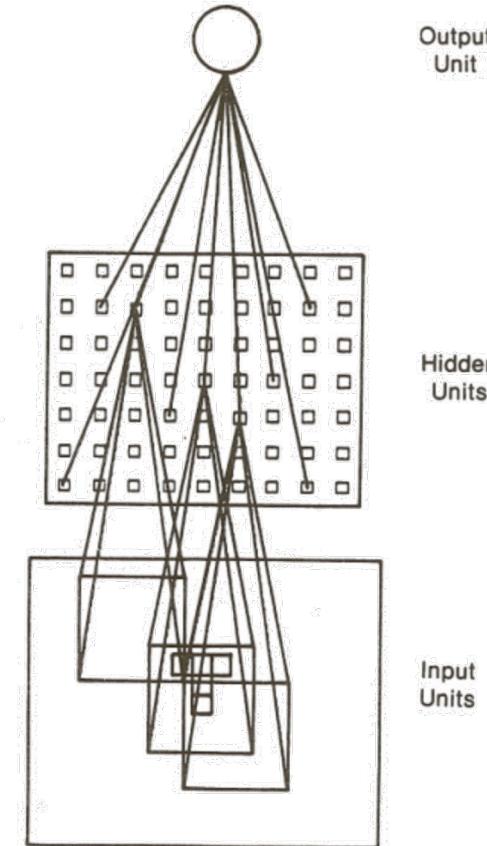
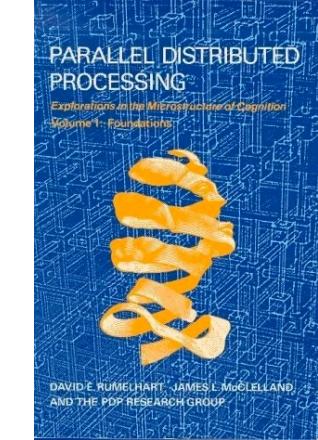


FIGURE 14. The network for solving the T-C problem. See text for explanation.



**PARALLEL DISTRIBUTED
PROCESSING**
**Explorations in the Microstructure
of Cognition**
Volume 1: Foundations

David E. Rumelhart James L. McClelland
and the PDP Research Group

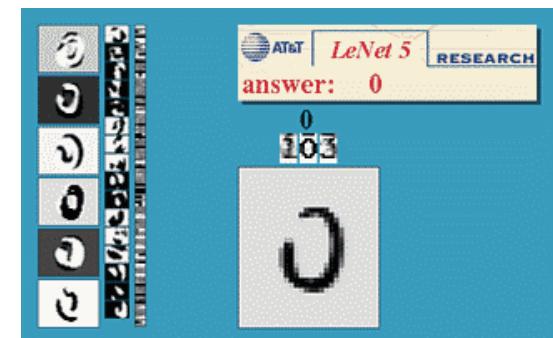
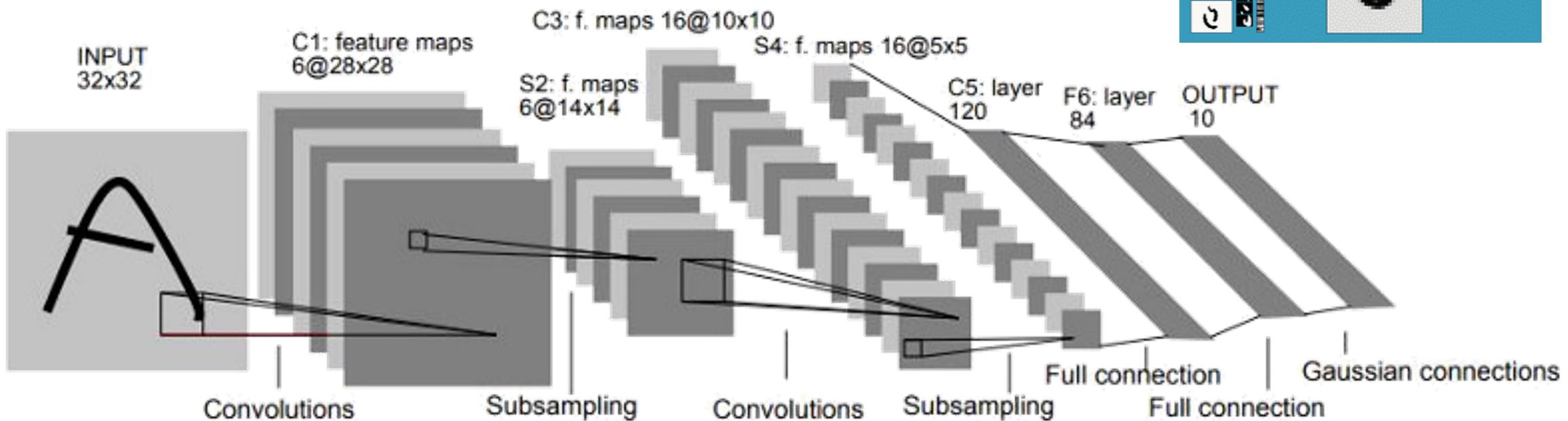
Chisato Asanuma Alan H. Kawamoto Paul Smolensky
Francis H. C. Crick Paul W. Munro Gregory O. Stone
Jeffrey L. Elman Donald A. Norman Ronald J. Williams
Geoffrey E. Hinton Daniel E. Rabin David Zipser
Michael I. Jordan Terrence J. Sejnowski

Institute for Cognitive Science
University of California, San Diego

1986

LeNet

[Backpropagation Applied to Handwritten Zip Code Recognition, LeCun et al., 1989]

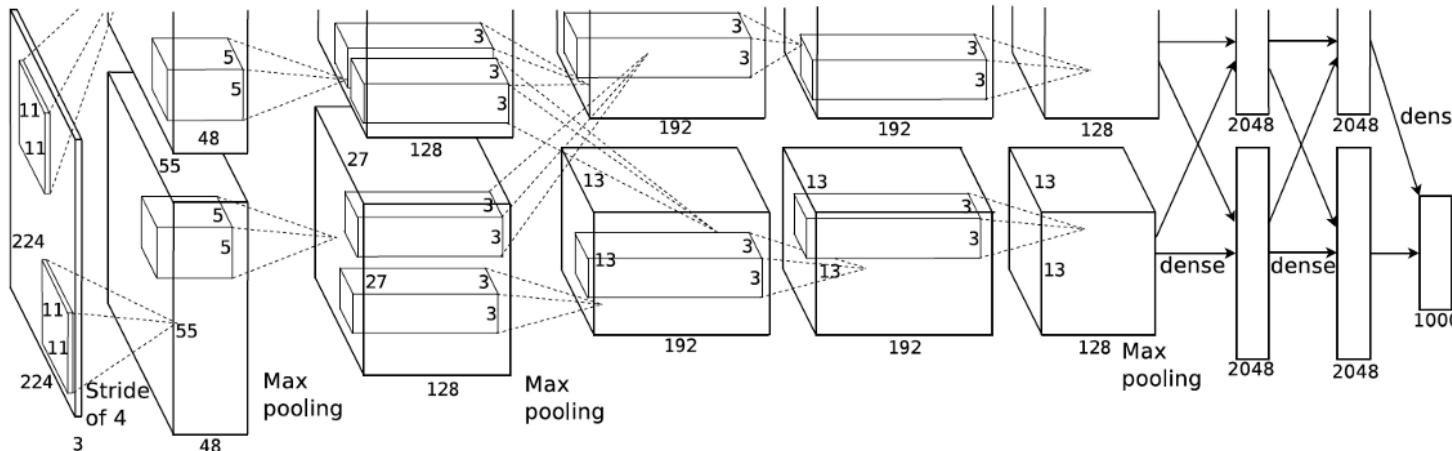


LeNet

1989

AlexNet

[Krizhevsky, Sutskever and Hinton, 2012]

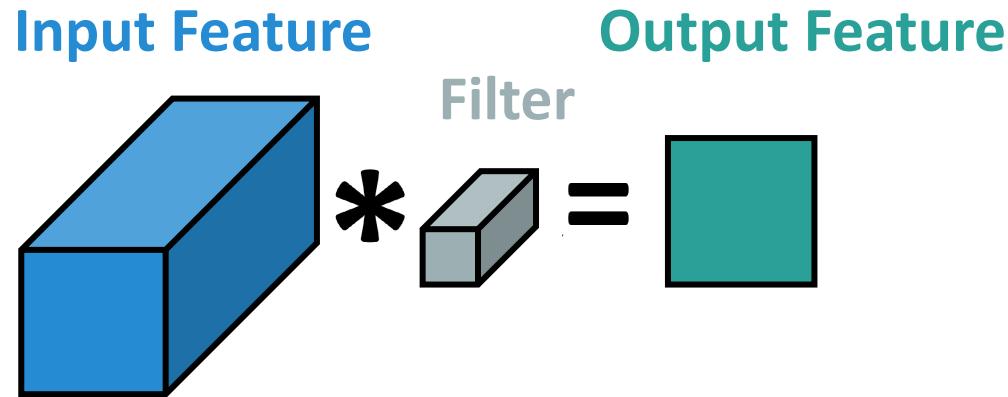


[Deng et al., 2009]
[Russakovsky et al., 2015]

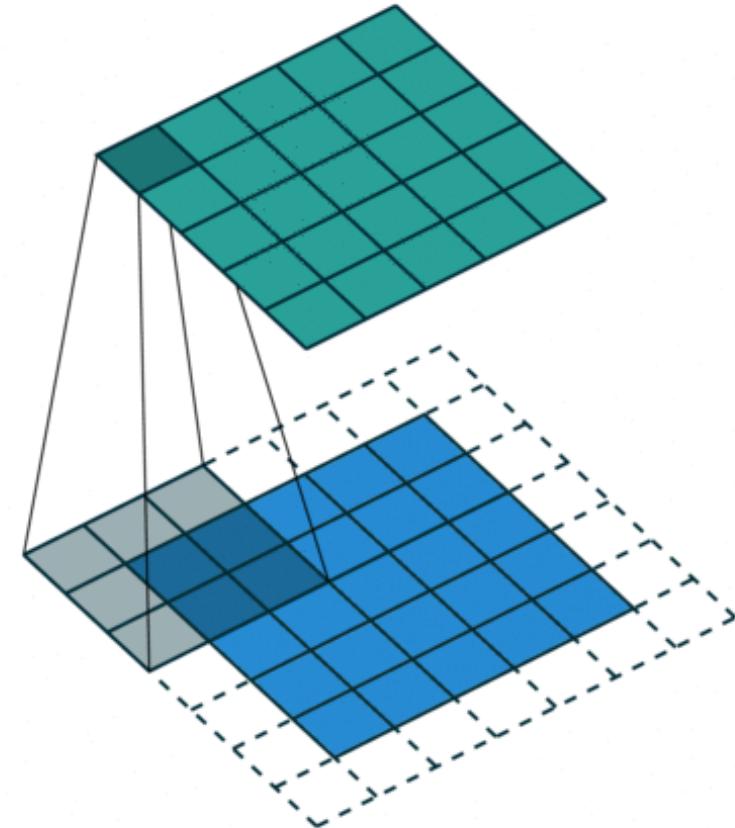
AlexNet

2012

ConvNet Basics: Convolution

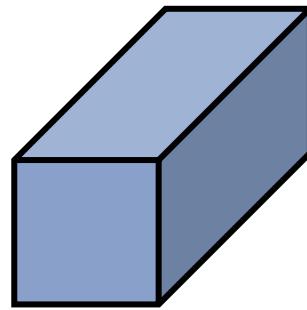


- Locally-connected
- Spatial weight-sharing
- Translation equivariance



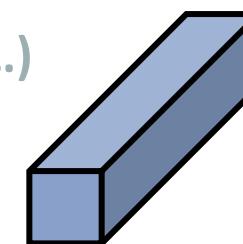
ConvNet Basics: Pooling

Input Feature



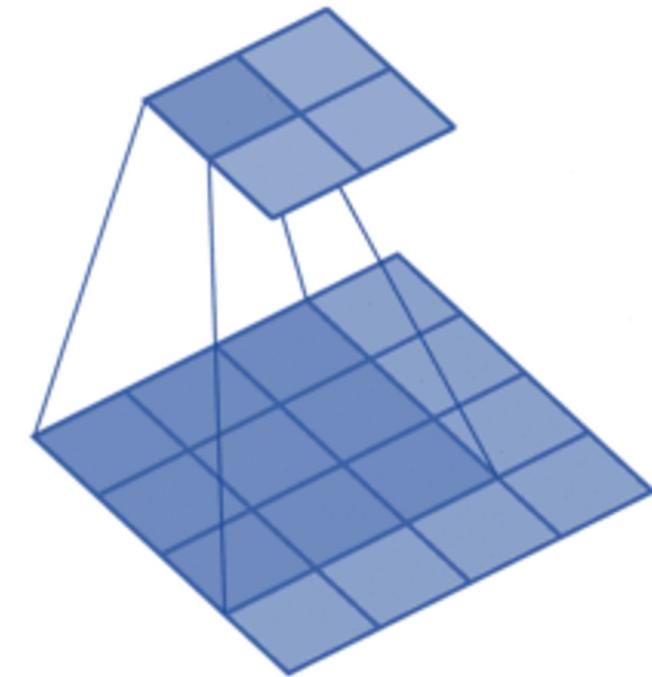
$W \times H \times C$

Output Feature



$W/2 \times H/2 \times C$

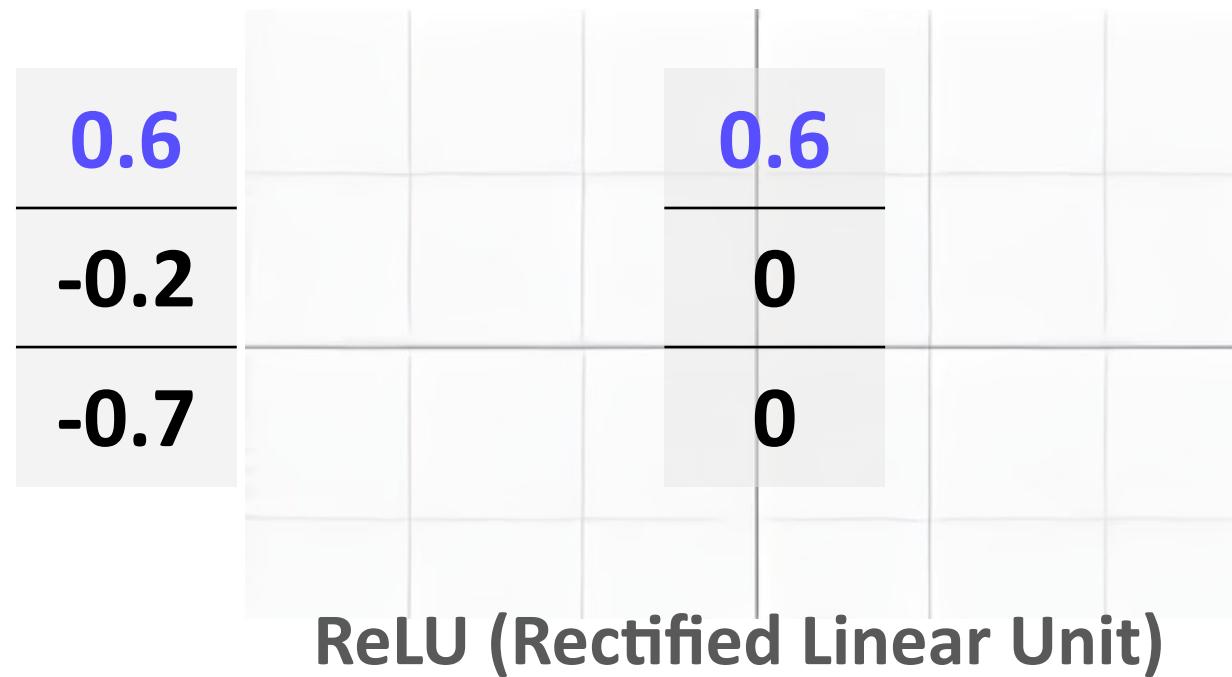
Pool (Max, Avg, ...)



- Invariant to small translations
- Reduce spatial resolutions



ConvNet Basics: Activation Function

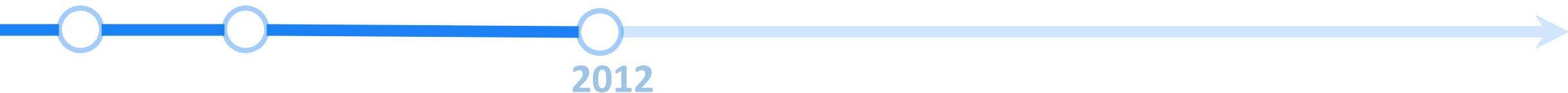
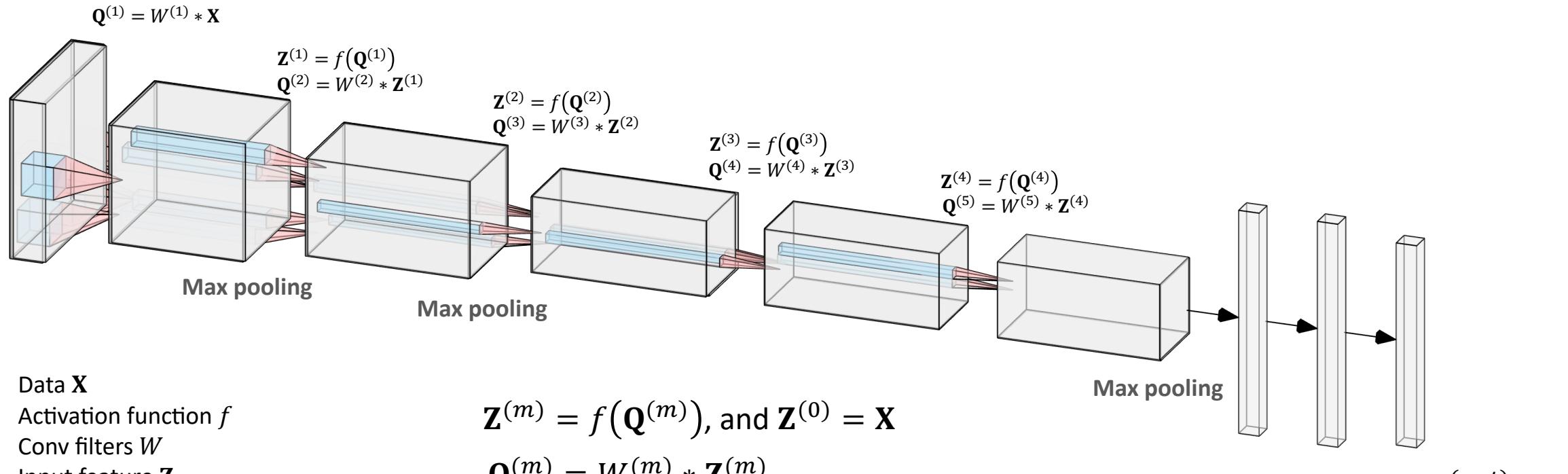


ConvNet Basics: Activation Function

Identity	Sigmoid	TanH	ArcTan
ReLU	Leaky ReLU	Randomized ReLU	Parameteric ReLU
Binary	Exponential Linear Unit	Soft Sign	Inverse Square Root Unit (ISRU)

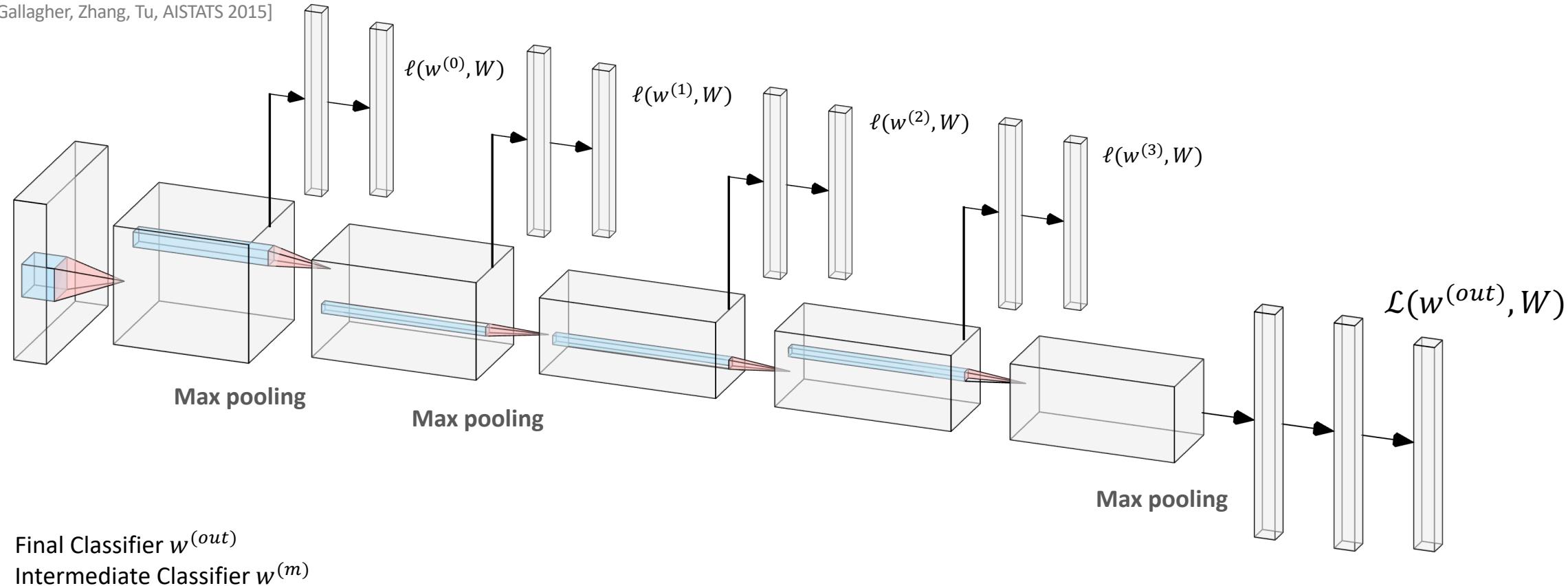


ConvNet Basics: Stacking layers



Deeply-supervised Nets

[Lee*, Xie*, Gallagher, Zhang, Tu, AISTATS 2015]



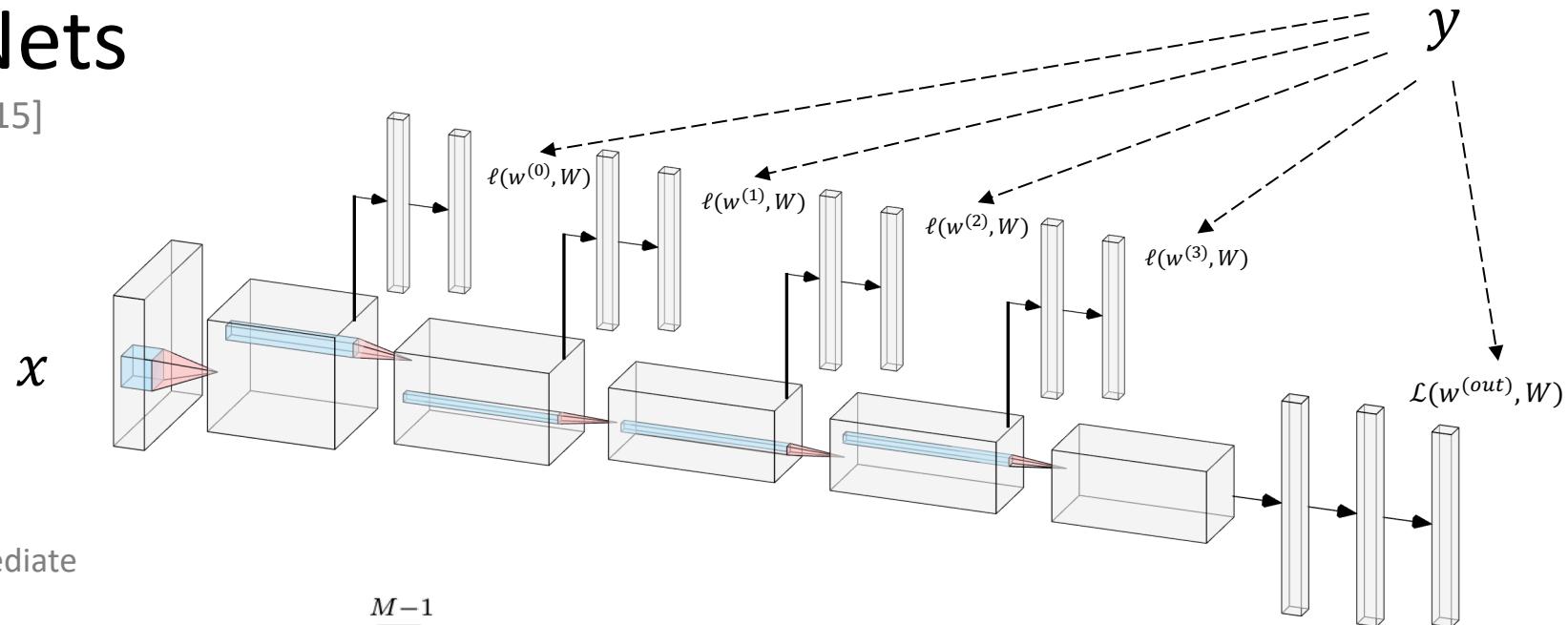
AlexNet DSN

2012

2014

Deeply-supervised Nets

[Lee*, Xie*, Gallagher, Zhang, Tu, AISTATS 2015]



Using hinge loss (SVM) for both output and intermediate classifiers

$$\text{Final training objective} \quad \|\mathbf{w}^{(out)}\|^2 + \mathcal{L}(W, \mathbf{w}^{(out)}) + \sum_{m=1}^{M-1} \alpha_m [\|\mathbf{w}^{(m)}\|^2 + \ell(W, \mathbf{w}^{(m)}) - \gamma]_+,$$

where

$$\mathcal{L}(W, \mathbf{w}^{(out)}) = \sum_{y_k \neq y} [1 - \langle \mathbf{w}^{(out)}, \phi(\mathbf{Z}^{(M)}, y) - \phi(\mathbf{Z}^{(M)}, y_k) \rangle]^2_+ \quad \text{Overall loss}$$

and

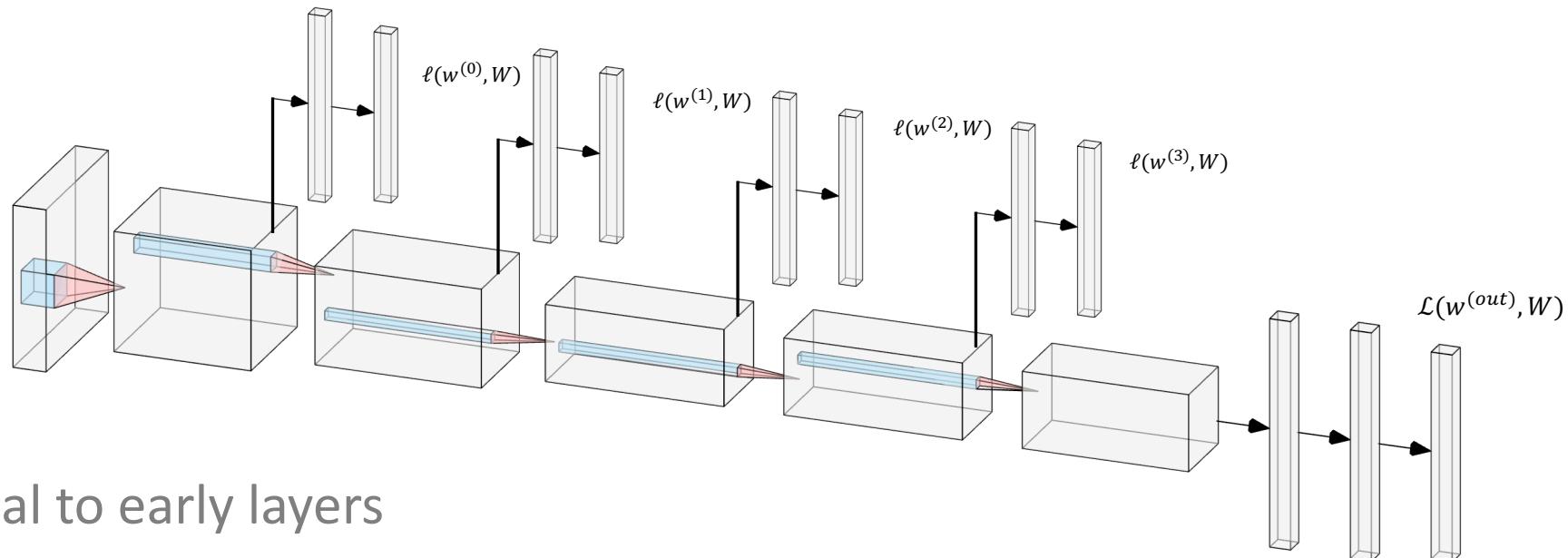
$$\ell(W, \mathbf{w}^{(m)}) = \sum_{y_k \neq y} [1 - \langle \mathbf{w}^{(m)}, \phi(\mathbf{Z}^{(m)}, y) - \phi(\mathbf{Z}^{(m)}, y_k) \rangle]^2_+ \quad \text{Companion losses}$$

AlexNet DSN



Deeply-supervised Nets

[Lee*, Xie*, Gallagher, Zhang, Tu, AISTATS 2015]

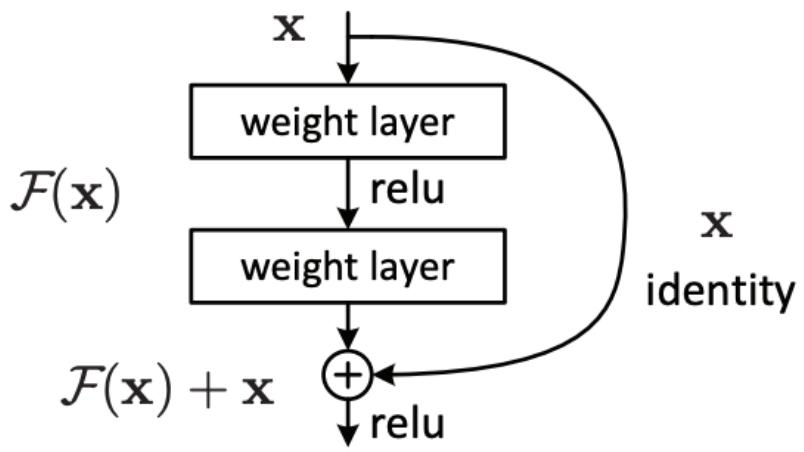


- Exposing training signal to early layers
 - Better regularization
 - Better convergence
- Popularized the idea of deep supervision (2000+ citations)

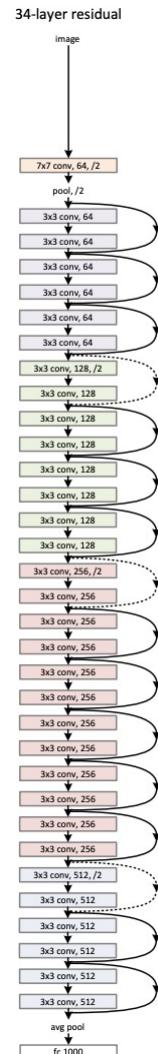


ResNet

[He et al., 2015]



repeating motif: a residual block



ResNet

2015



figure: reddit?

ResNeXt

[Xie, Girshick, Dollár, Tu, He, CVPR 2017]

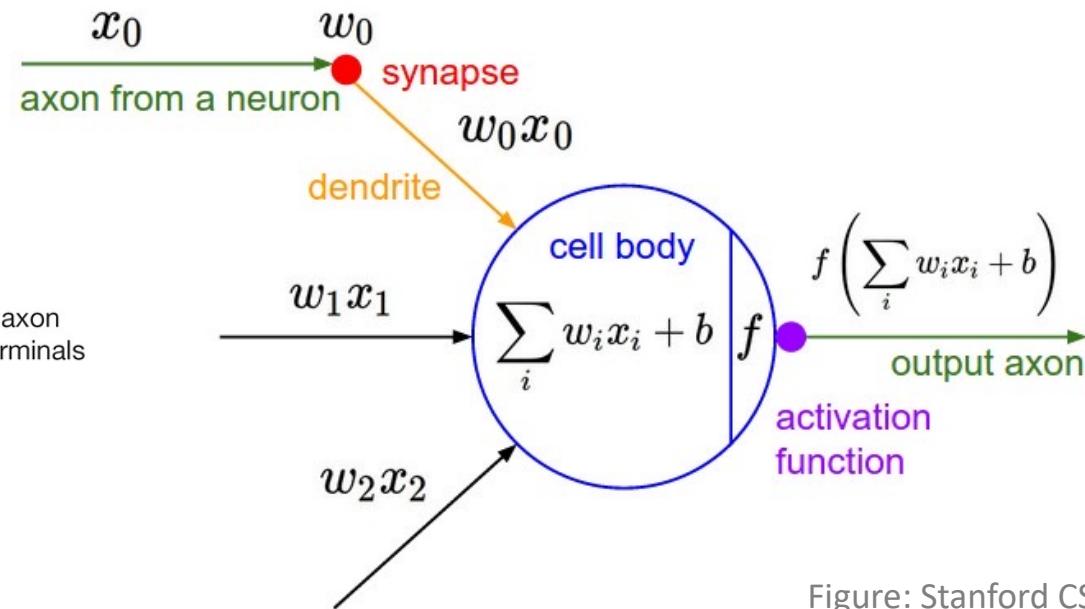
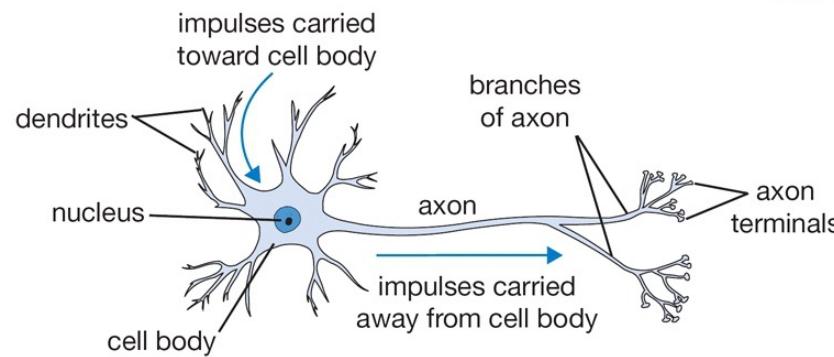


Figure: Stanford CS231n

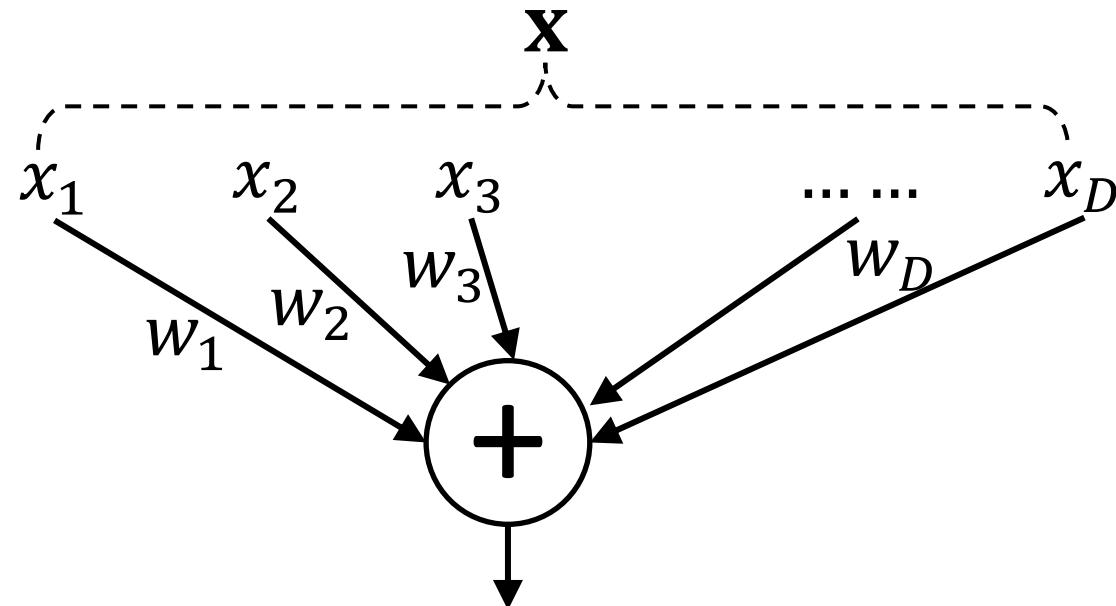


ResNeXt

[Xie, Girshick, Dollár, Tu, He, CVPR 2017]

Analysis of a simple neuron:

- Splitting: $\mathbf{x} \rightarrow x_i$
- Transforming: $w_i x_i$
- Aggregating: $\sum_{i=1}^D w_i x_i$

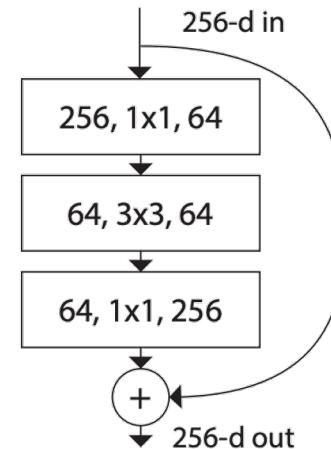


ResNeXt

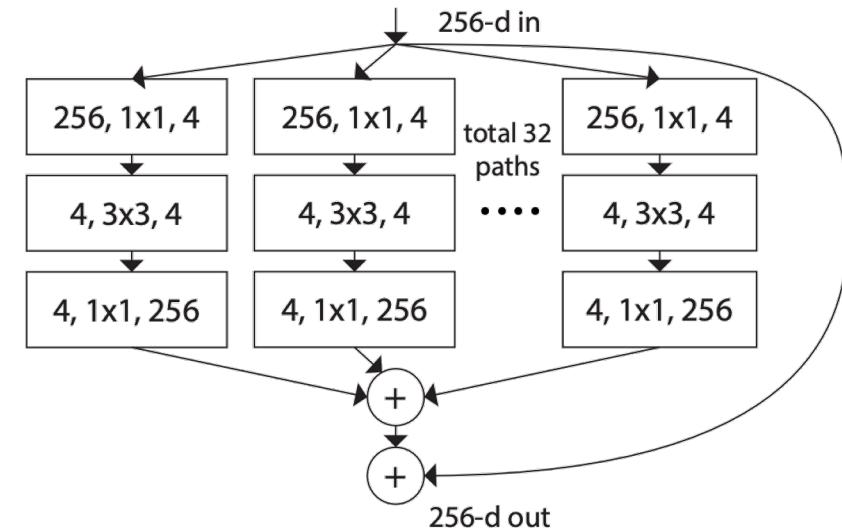
[Xie, Girshick, Dollár, Tu, He, CVPR 2017]

“Network-in-Neuron”:

- Splitting: Parallel pathways
- Transforming: Conv block
- Aggregating: $\mathcal{F}(\mathbf{x}) = \sum_{i=1}^C \mathcal{T}_i(\mathbf{x})$



ResNet: Single Stream

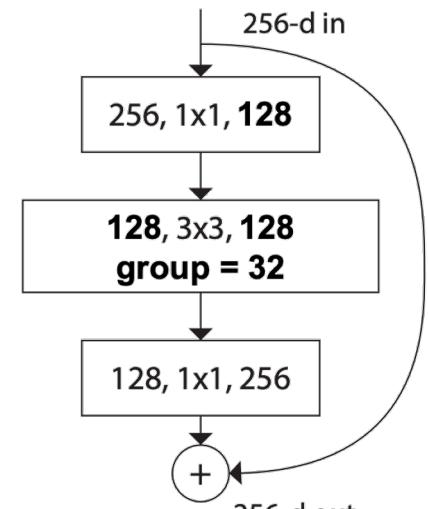
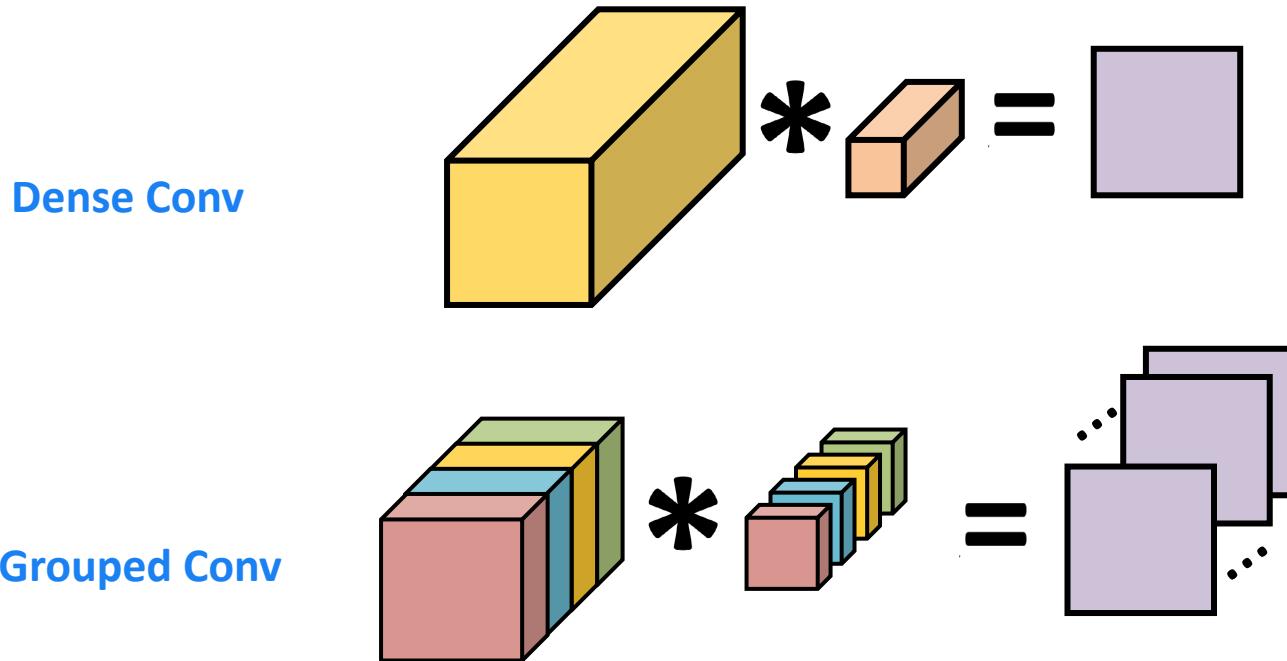


ResNeXt: Multiple pathways



ResNeXt

[Xie, Girshick, Dollár, Tu, He, CVPR 2017]



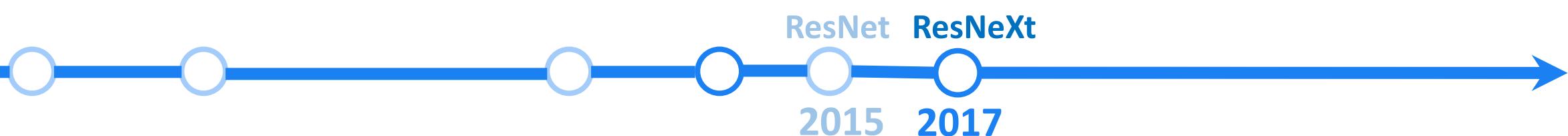
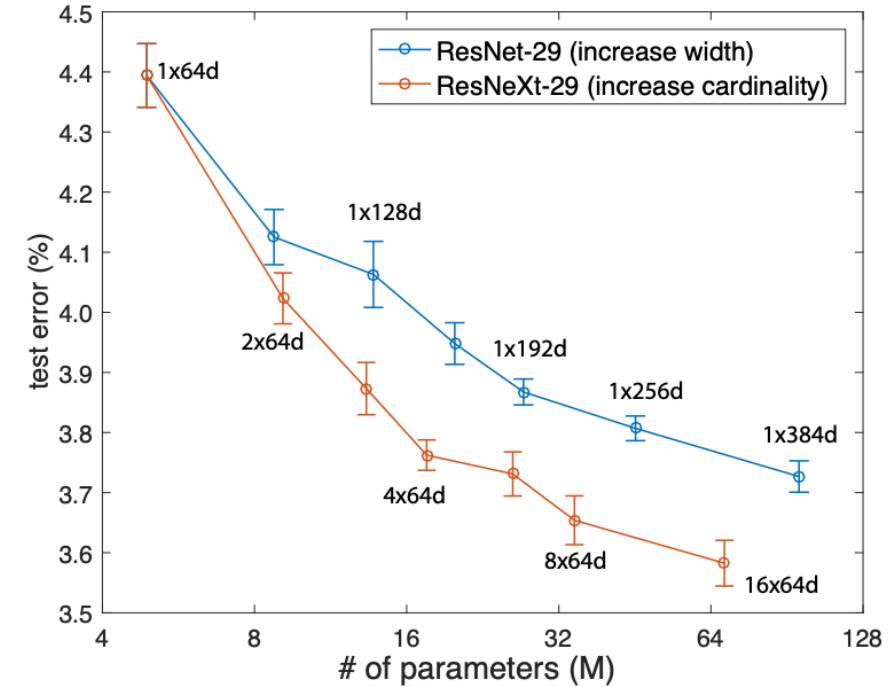
a ResNeXt block

ResNeXt
2016

ResNeXt

[Xie, Girshick, Dollár, Tu, He, CVPR 2017]

- Scaling neural networks:
 - Going deeper / wider
 - We identify the **NeXt** dimension:
Cardinality (# number of groups)
- Increasing cardinality:
 - Better capacity-accuracy trade-off
 - Wider and sparser
 - Implicit ensemble
 - Disentangled features

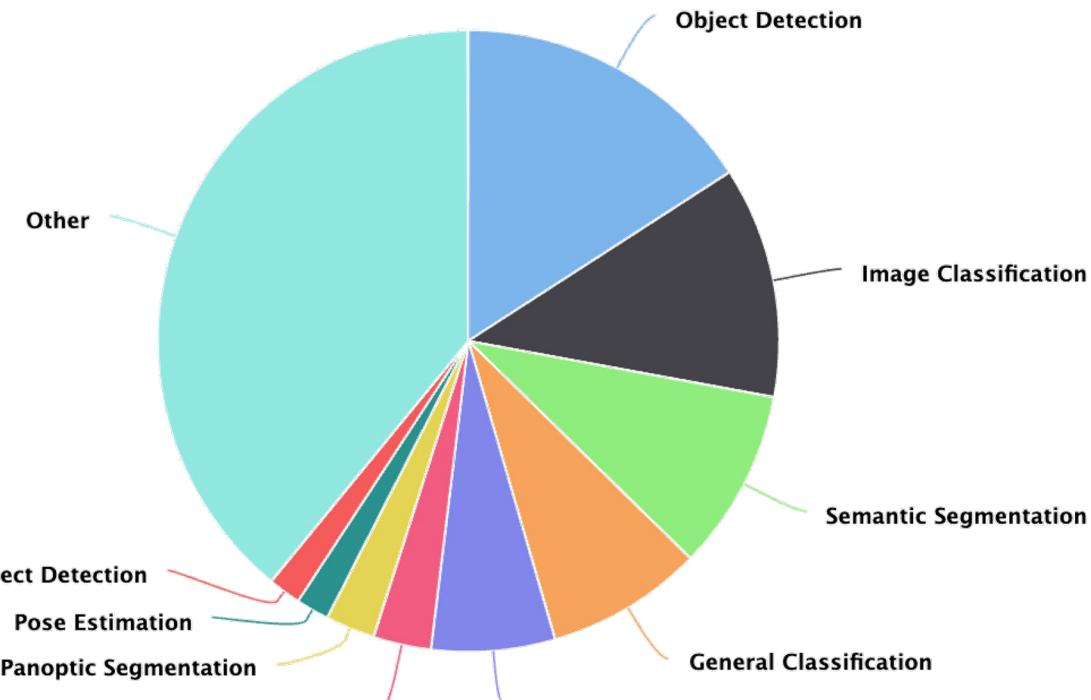


ResNeXt

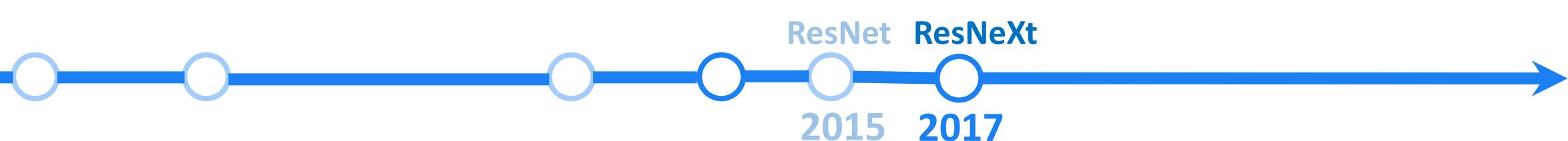
[Xie, Girshick, Dollár, Tu, He, CVPR 2017]

- Better feature learning
- One of the most popular visual backbones
- Widely used in various academic and industrial applications

ResNeXt empowers many visual recognition tasks



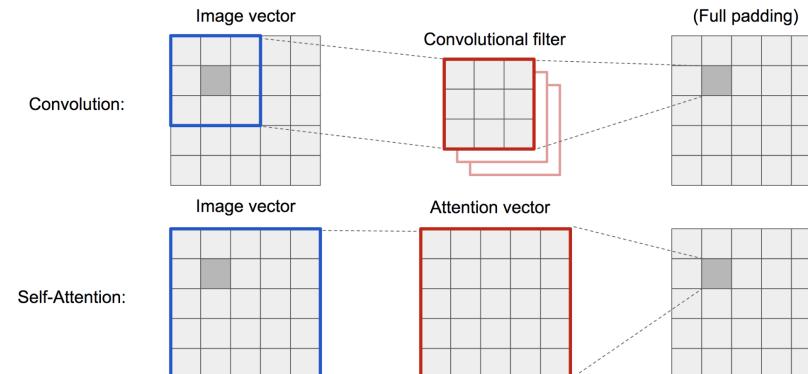
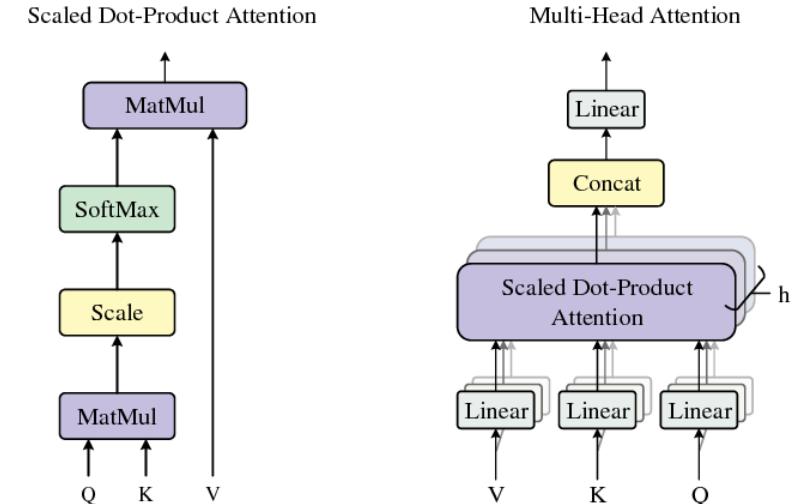
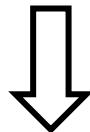
Source: paperswithcode.com/method/resnext



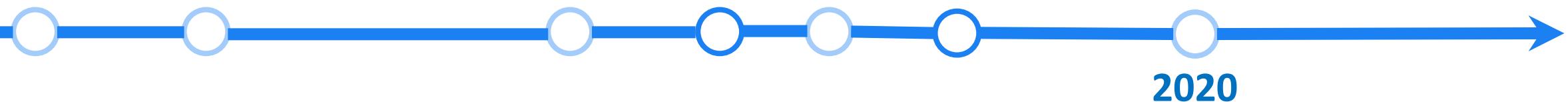
Vision Transformer (ViT)

[Dosovitskiy et al., ICLR 2021]

- Self-attention:
 - Less inductive bias
- Scalable
 - w/ bigger model
 - w/ larger data



courtesy: Lilian Weng



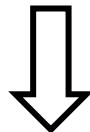
Vision is NOT just about classification

- Expanding vision capabilities
 - Large resolution
 - Multi-scale
- Vision Transformer
 - Quadratic complexity
 - No hierarchy



Advanced Vision Transformers

- Self-attention:
 - Less inductive bias



- Scalable
 - w/ bigger model
 - w/ larger data

Self-attention (ViT)

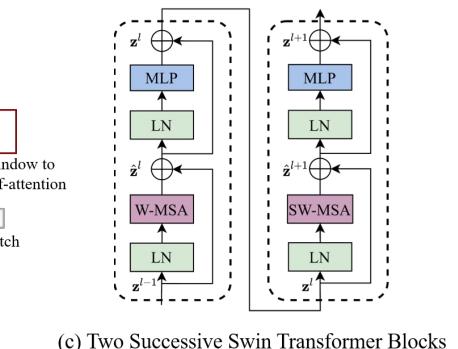
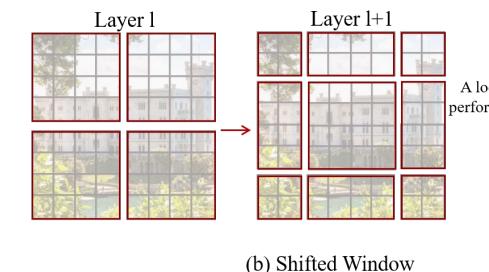
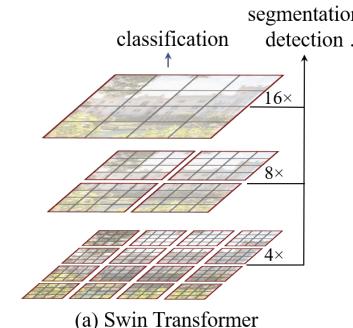
+ Conv Priors

Translation equivariance

Locality

Hierarchical

More complicated designs + Specialized modules



Swin Transformer

2021

ConvNet losing steam?

Venue	Convolution, CNN, ConvNet	Attention, “-Former”
ECCV 2020	56	54
CVPR 2021	49	78
ICCV 2021	44	176
CVPR 2022	44	263

Swin Transformer

 State of the Art Object Detection on COCO test-dev (using additional training data)

 State of the Art Instance Segmentation on COCO test-dev

 State of the Art Semantic Segmentation on ADE20K (using additional training data)

 Ranked #4 Action Classification on Kinetics-400 (using additional training data)

Swin Transformers is a *hybrid* architecture

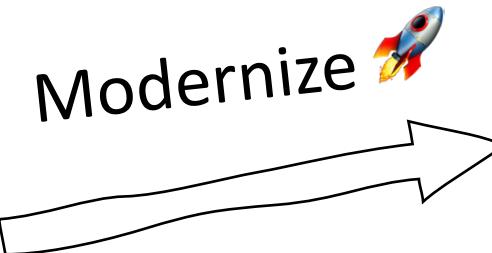
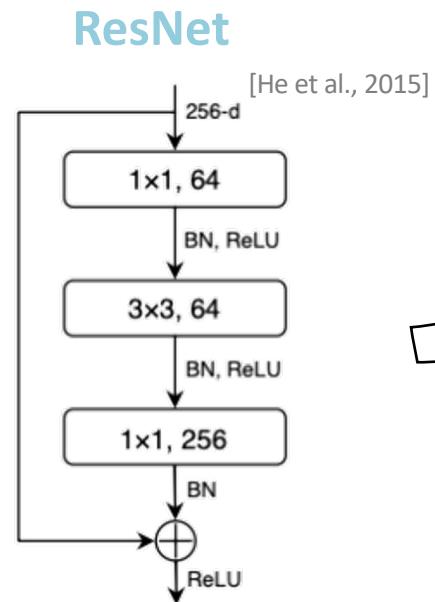
- Similarity:
 - Convolution inductive bias
- Difference:
 - “Core” component (attention vs. convolution)
 - Training procedures
 - Macro and micro architecture design decisions
- Common **assumption** in the 2020s:
 - **Self-attention** is the key for superior performance and scalability.
 - ConvNet is NOT a scalable architecture

A ConvNet for the 2020s

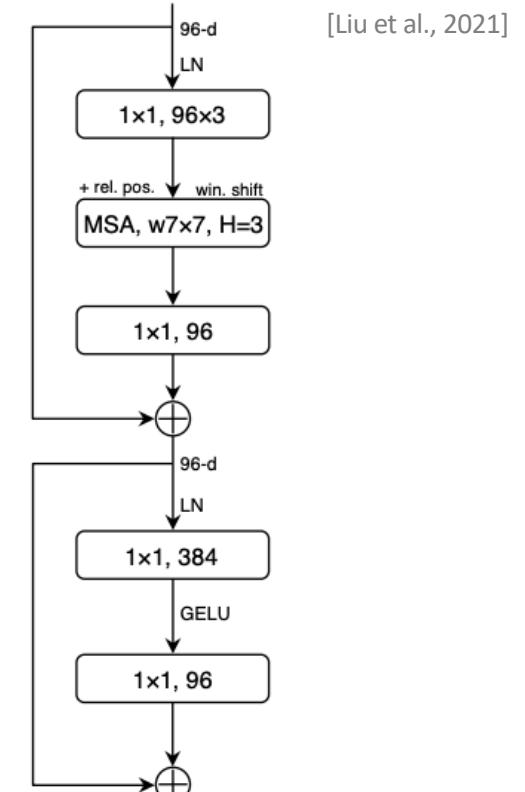
[Liu, Mao, Wu, Feichtenhofer, Darrell, Xie. CVPR 2022]

Central question:

- How do the design choices in **Transformers** impact a **ConvNet's** performance?



Hierarchical Vision Transformer



ResNet

2015

ResNeXt

2017

ViT

2020

ConvNeXt

2022

Improved training recipe

("DeiT Recipe") [Touvron et al., 2021]

Typical Vision Transformer Training Recipe



Typical ResNet Training Recipe



ResNet-50 ImageNet top-1: 76.7% -> 78.8% 🚀

[Revisiting ResNets: Improved Training and Scaling Strategies, Bello et al, 2021]

[ResNet strikes back: An improved training procedure in timm Wightman, et al, 2021]

Swin-T/B

ImageNet
Top1 Acc (%)

78

80

82

4.5

ResNet-50/200

78.8

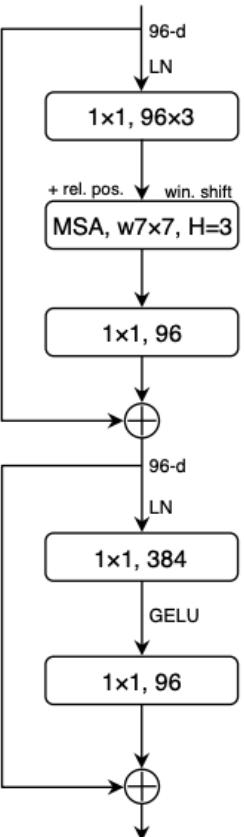
GFLOPs
4.1

Macro Design

- stage ratio
- “patchify” stem

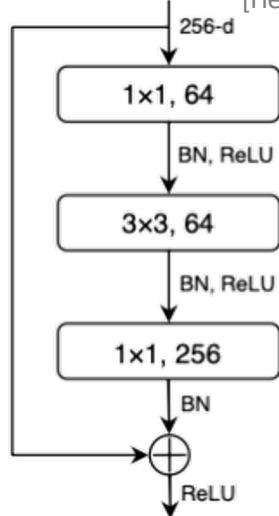
SwinTransformer

[Liu et al., 2021]



ResNet

[He et al., 2015]



Modernize

Swin-T/B

ImageNet
Top1 Acc (%)

78

81.3

4.5

80

82

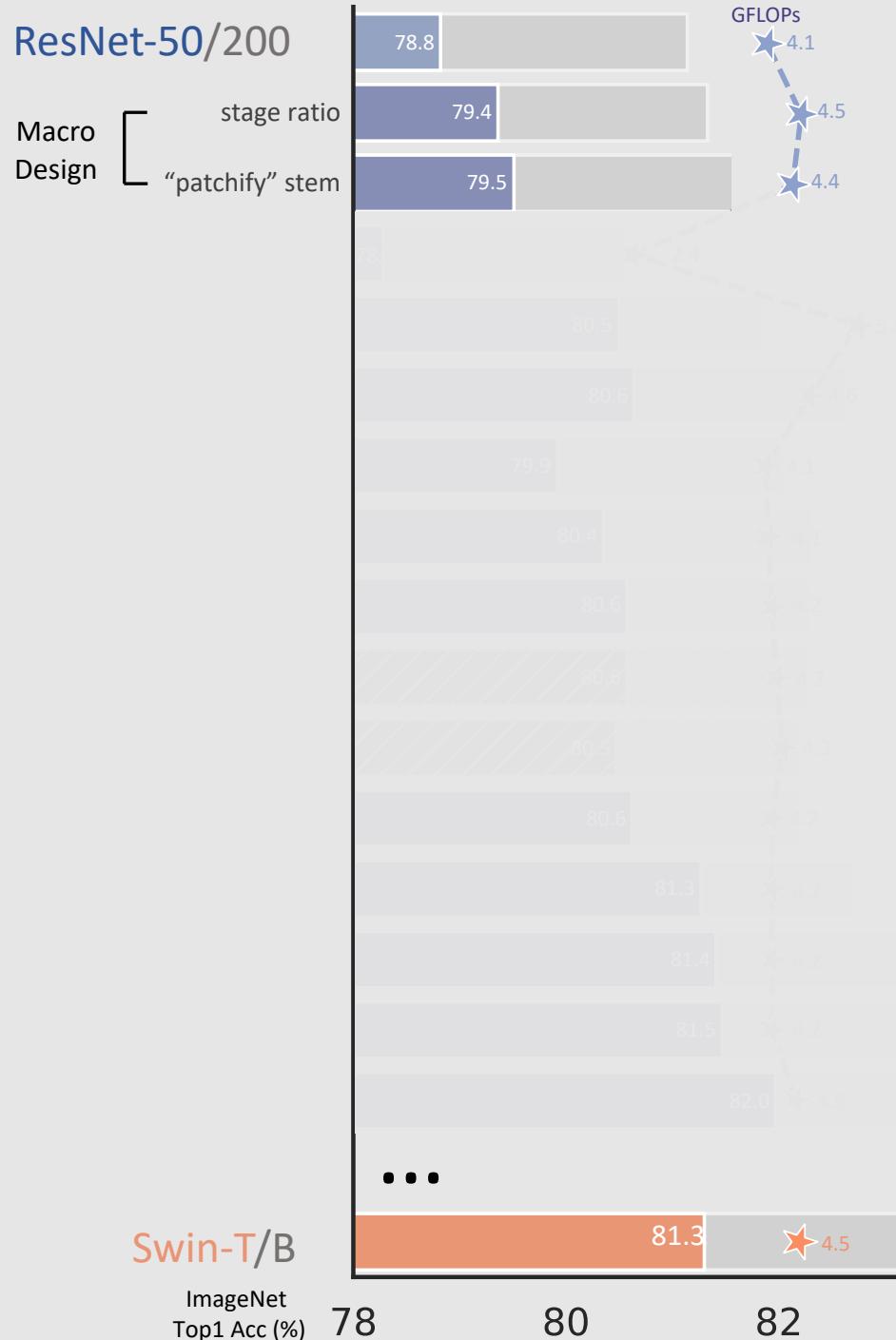
78 80 82

Input Stem (ResNet)

layer name	output size	18-layer	34-layer	50-layer	101-layer	152-layer
conv1	112×112			7×7, 64, stride 2		
				3×3 max pool, stride 2		
conv2_x	56×56	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$
conv3_x	28×28	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 8$
conv4_x	14×14	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 23$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 36$
conv5_x	7×7	$\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$
	1×1	average pool, 1000-d fc, softmax				
FLOPs		1.8×10^9	3.6×10^9	3.8×10^9	7.6×10^9	11.3×10^9

ResNet-50/200

Macro Design [stage ratio
“patchify” stem]



Swin-T/B

ImageNet
Top1 Acc (%)

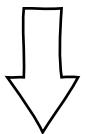
78

80

82

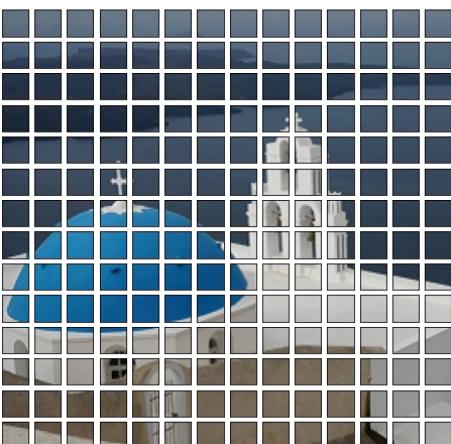
Input Stem

Overlapping Conv + Max Pooling



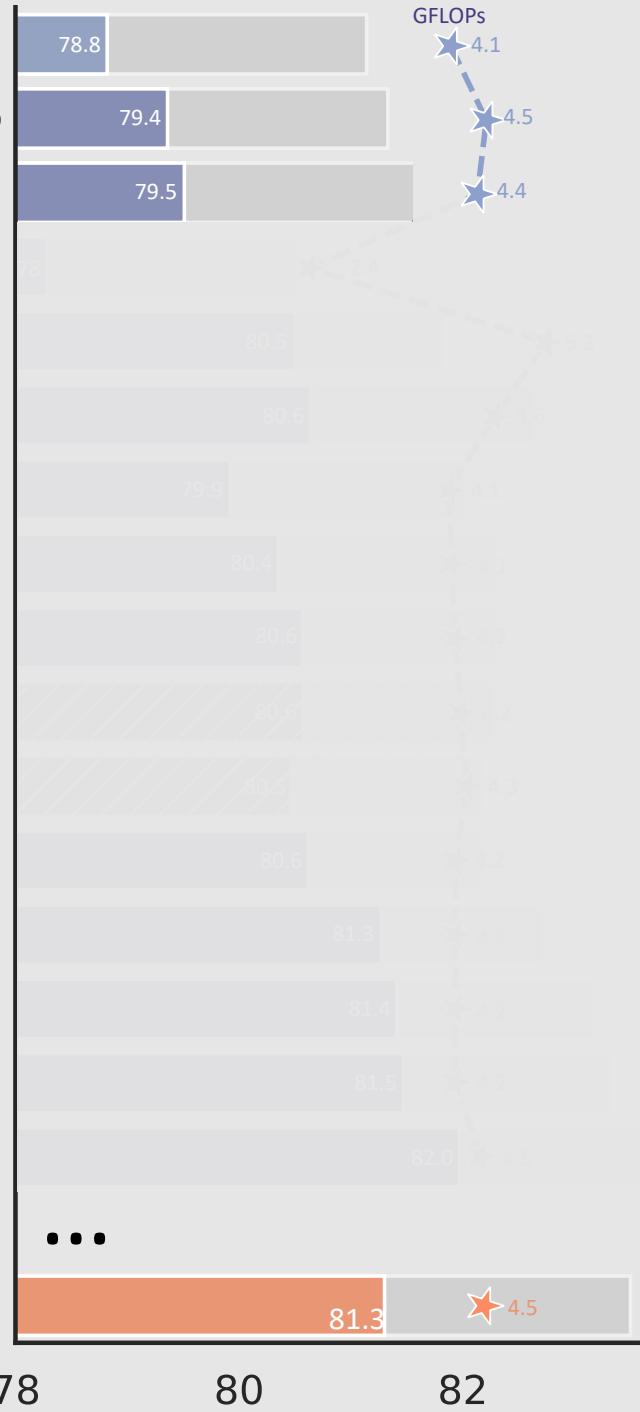
Non-overlapping Conv (4x4, stride 4)

(a.k.a Patchify)



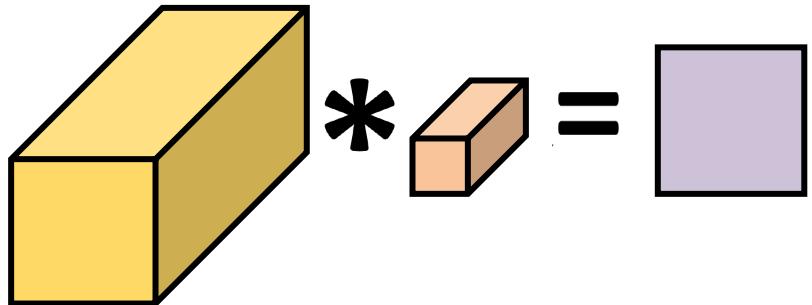
ResNet-50/200

Macro Design [stage ratio
"patchify" stem

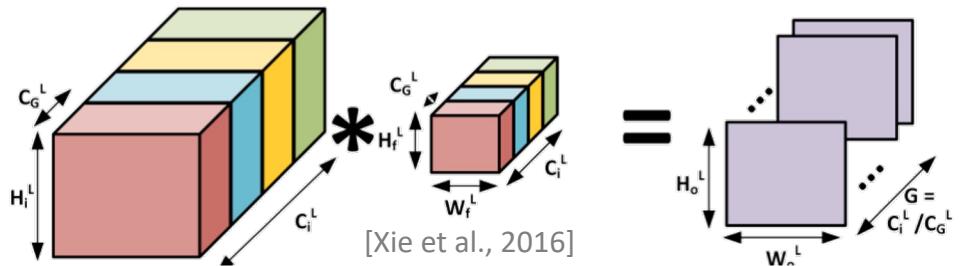


ResNeXt-ify

Dense Conv

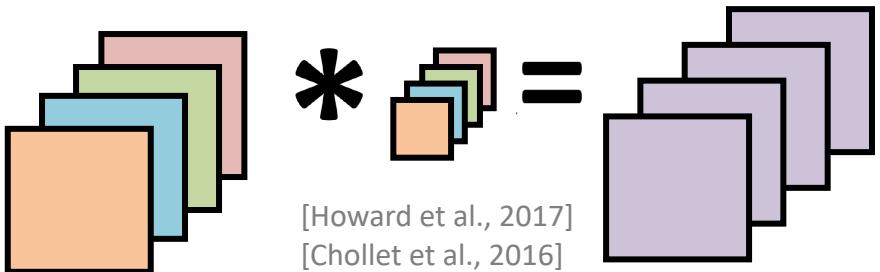


Grouped Conv

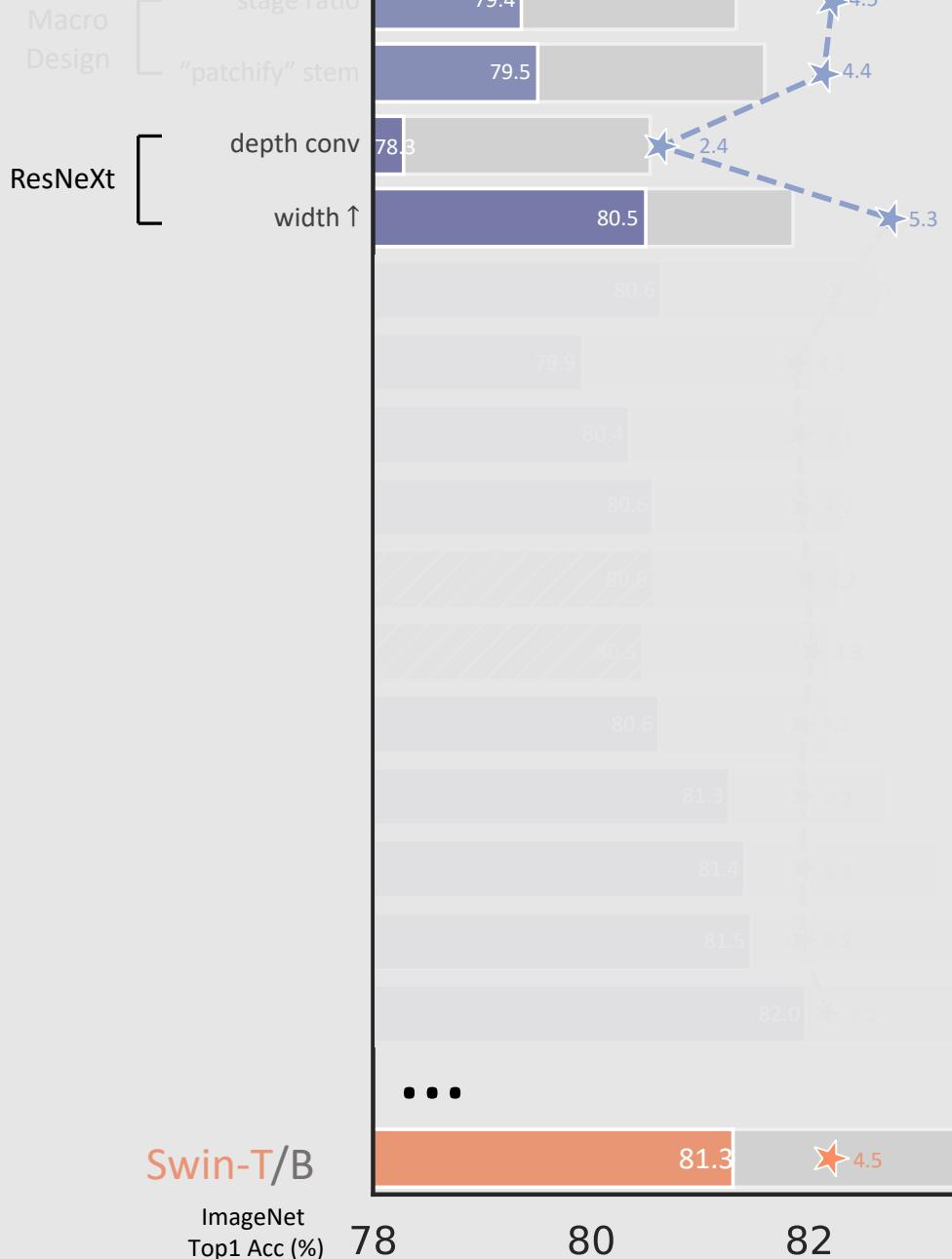


Depthwise Conv

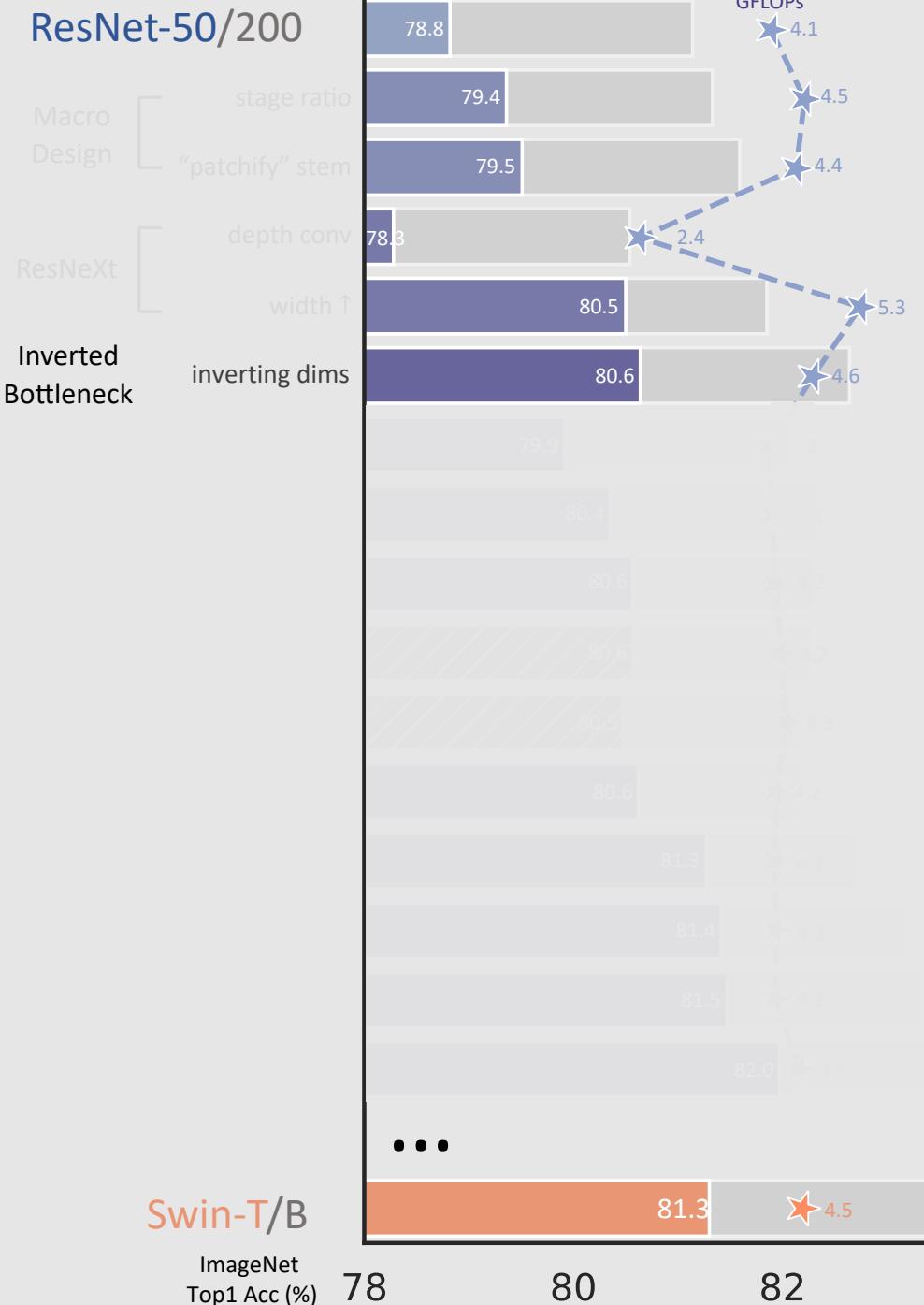
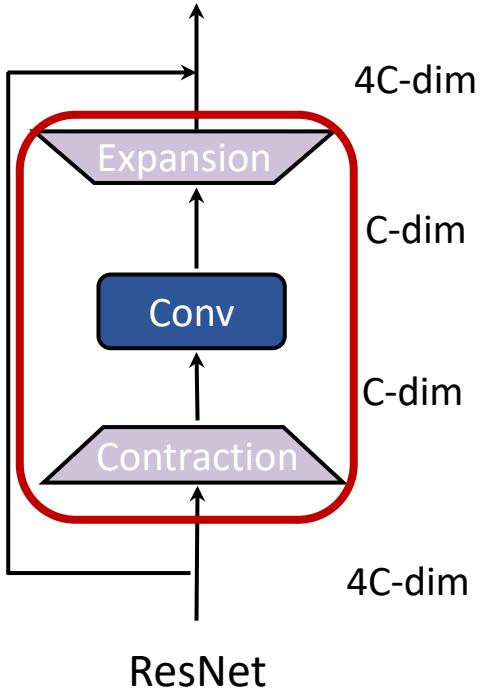
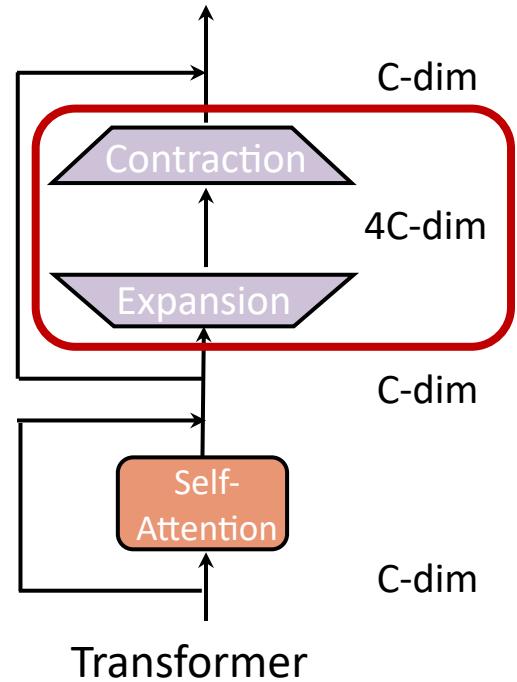
groups = # channels



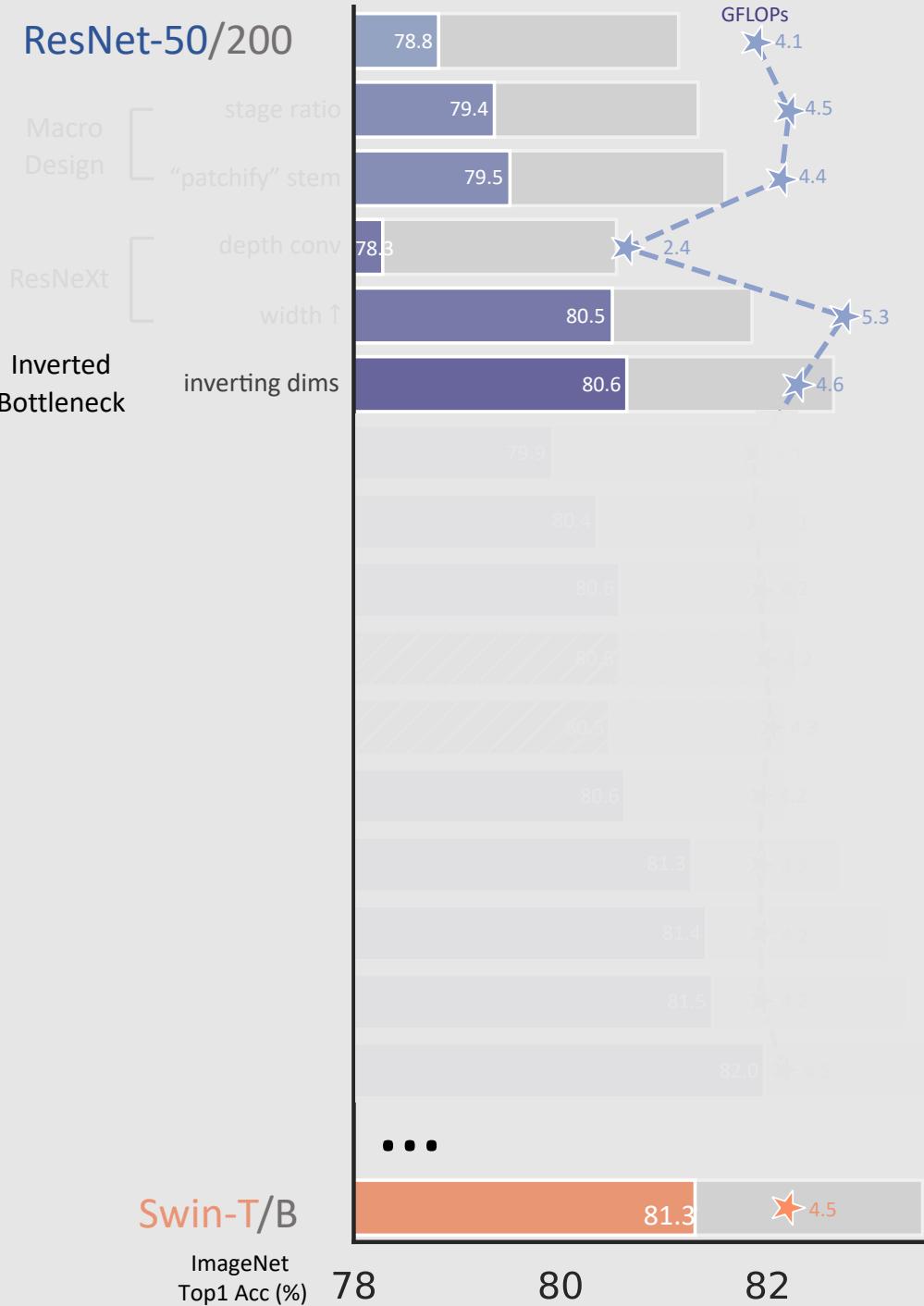
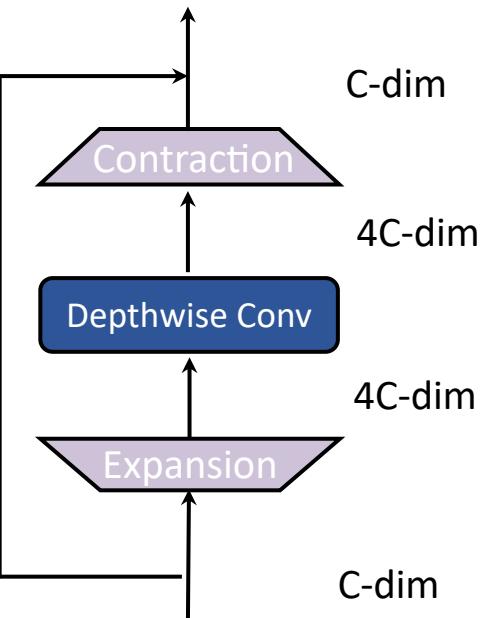
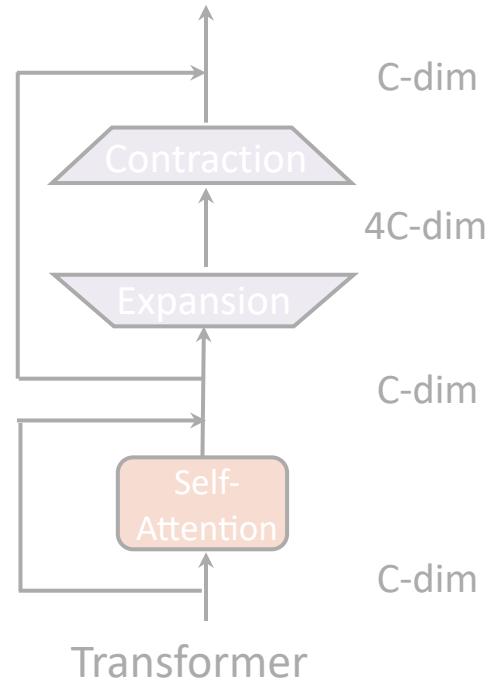
ResNet-50/200



Inverted Bottleneck

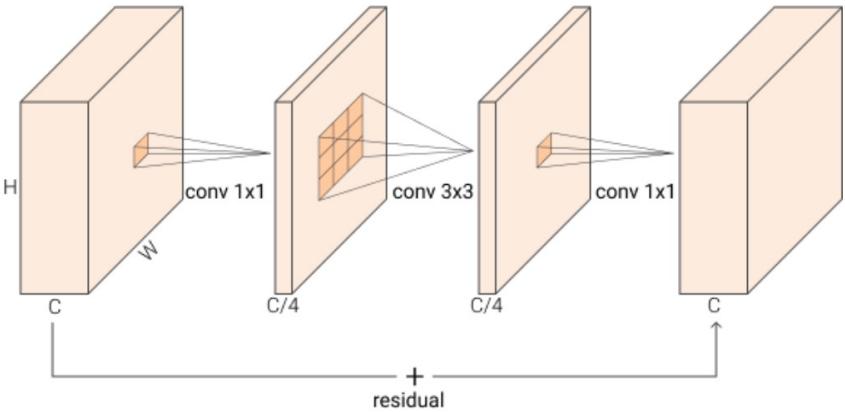


Inverted Bottleneck



2016

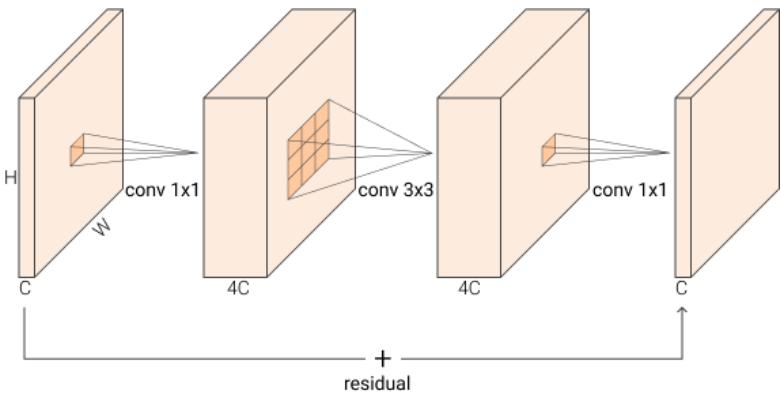
ResNet



2018

Inverted Bottleneck MobileNet v2

[Sandler et al, 2018]



ResNet-50/200

Macro Design

- "patchify" stem
- depth conv
- width ↑
- inverting dims

ResNeXt

stage ratio

- 78.8
- 79.4
- 79.5
- 78.3
- 80.5
- 80.6
- 79.9
- 80.4
- 80.6
- 80.8
- 80.5
- 80.6
- 81.3
- 81.4
- 81.5
- 82.0

Inverted Bottleneck

GFLOPs

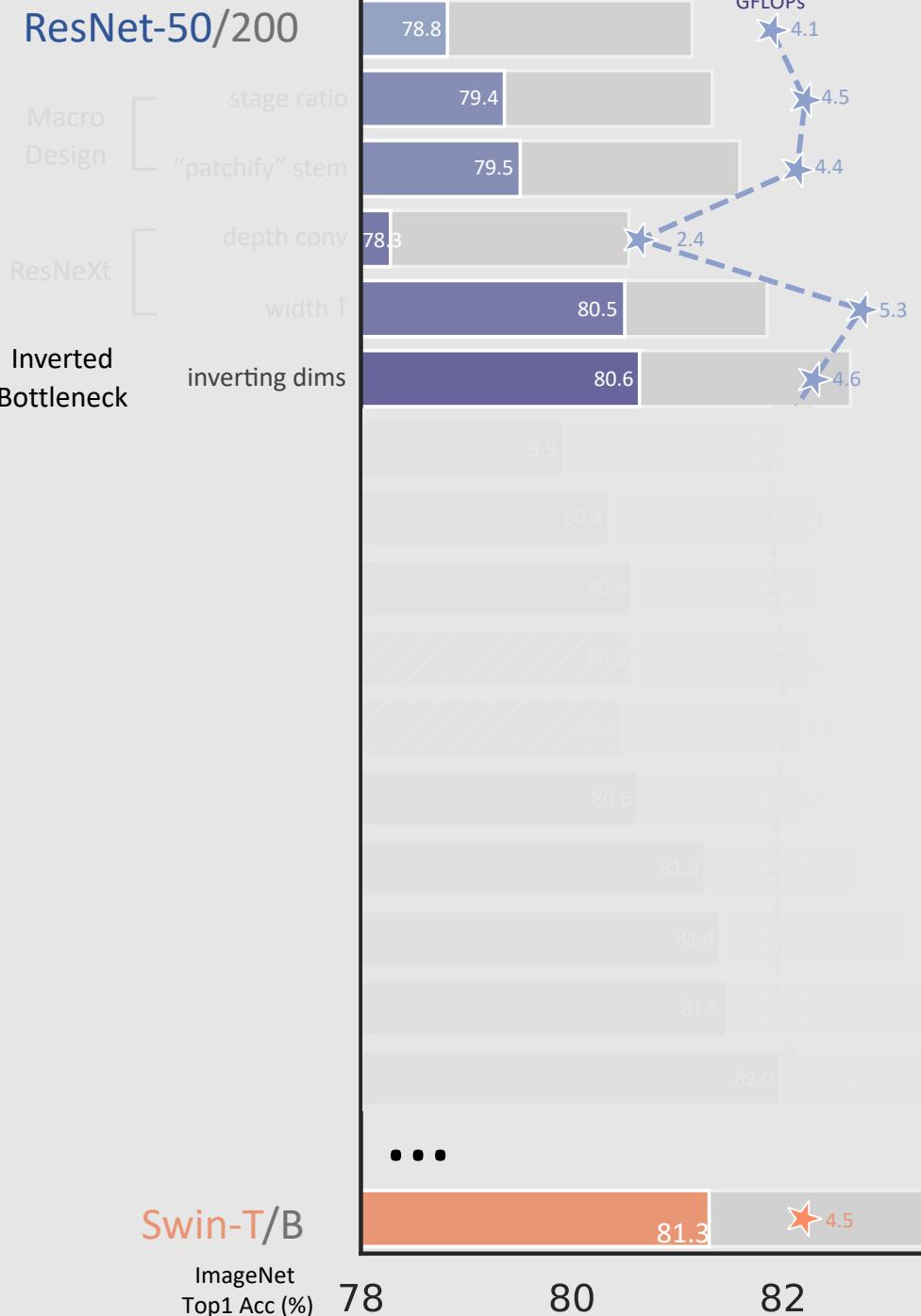
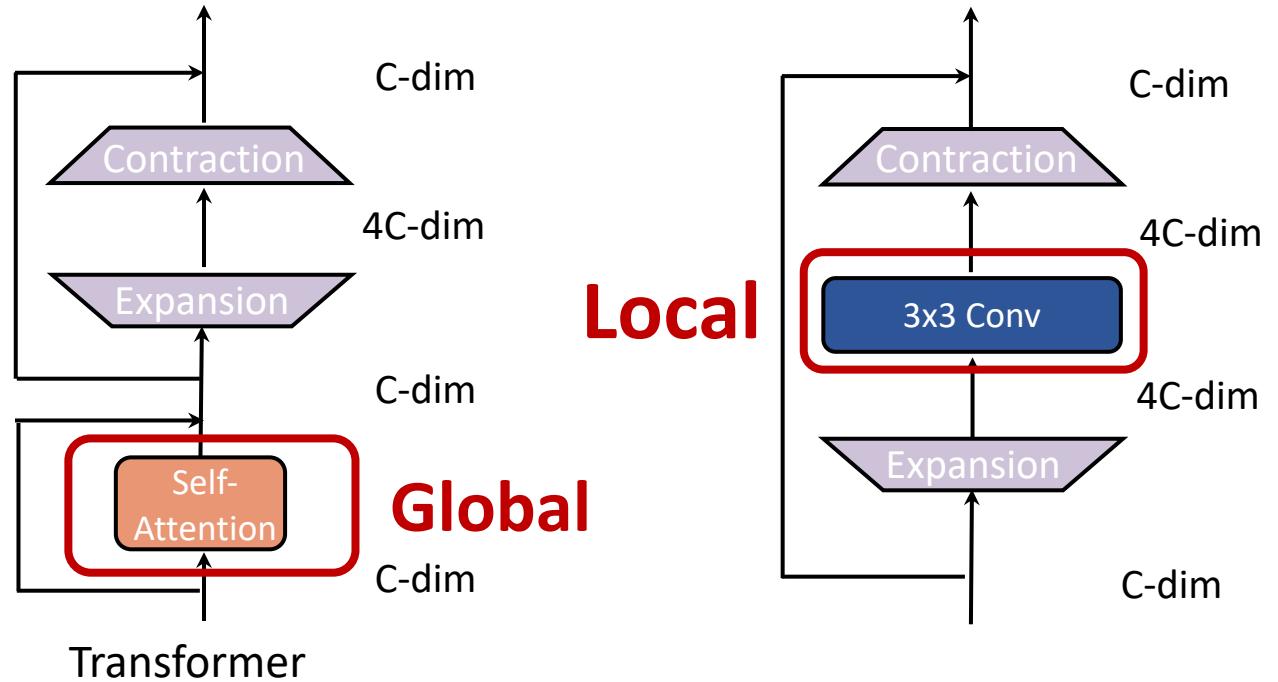
- 4.1
- 4.5
- 2.4
- 5.3
- 4.6

Swin-T/B

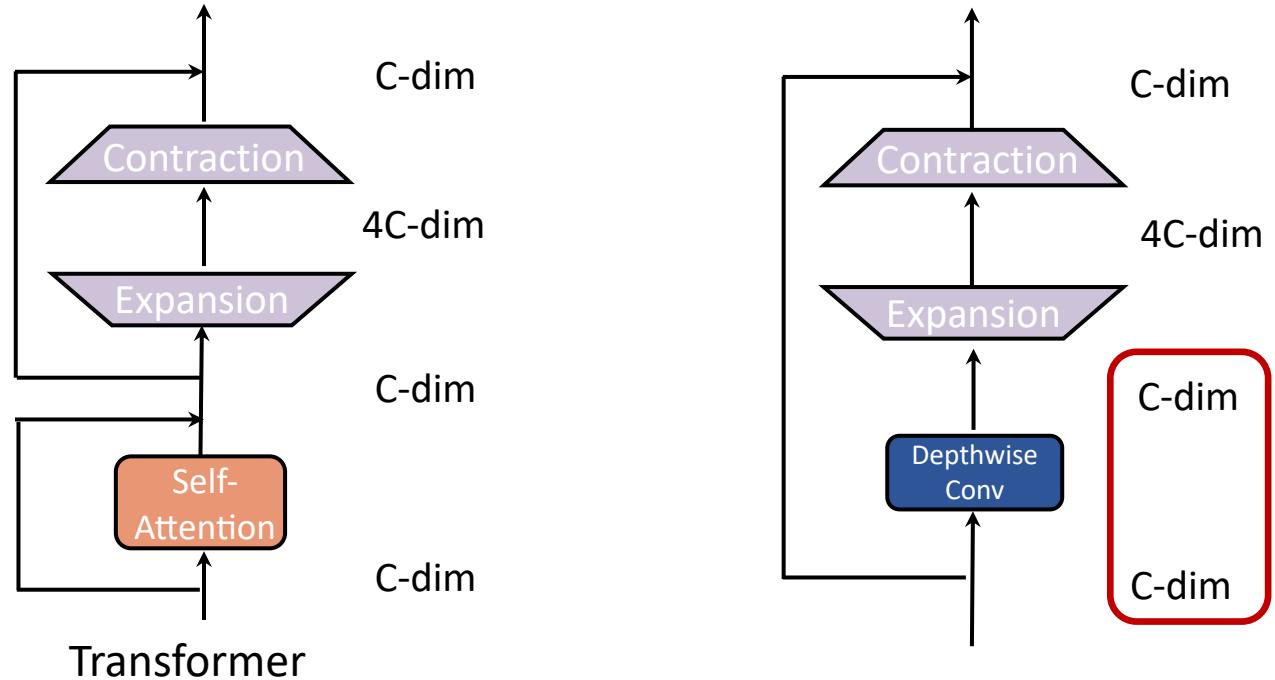
ImageNet
Top1 Acc (%)



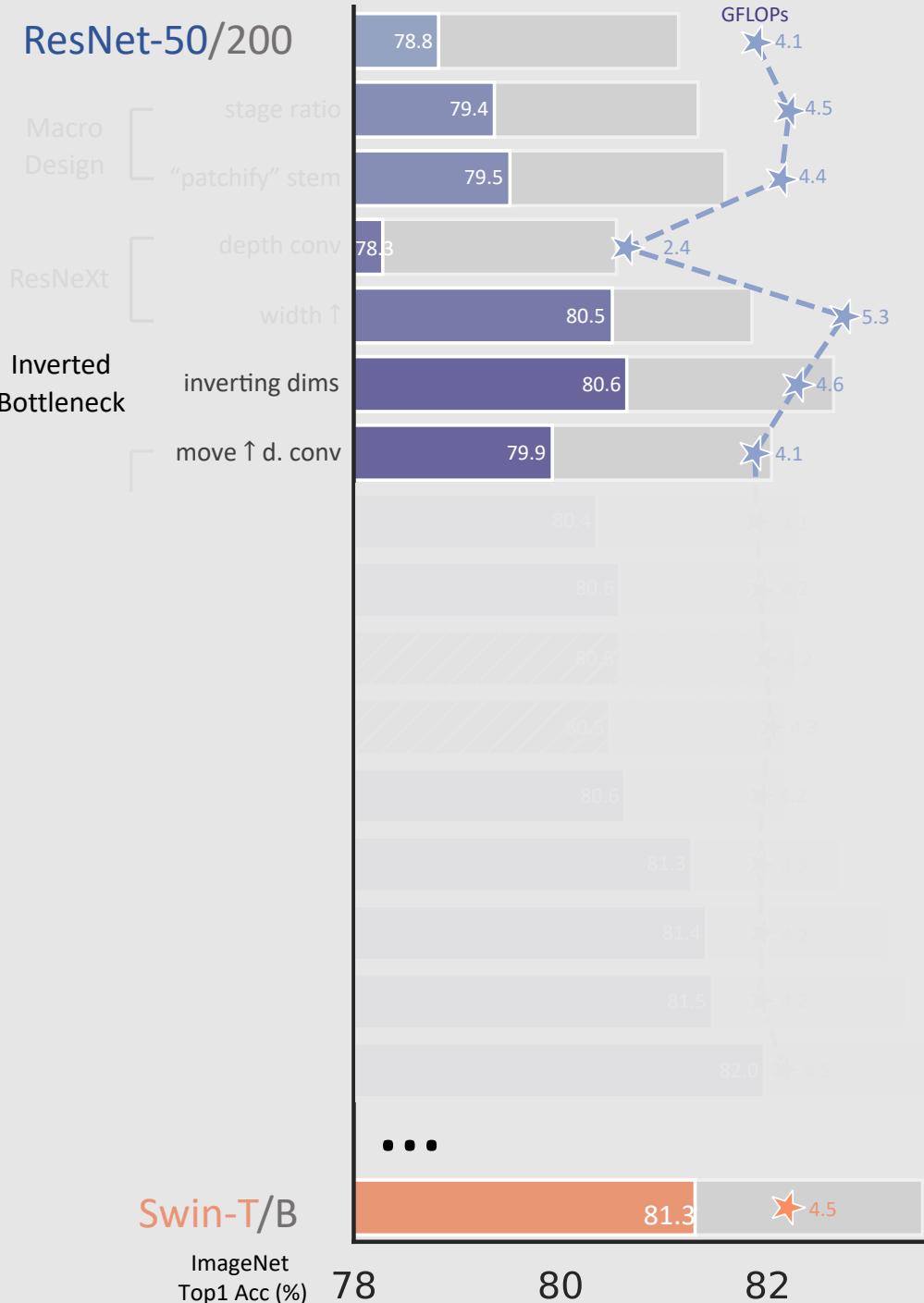
Large kernel size



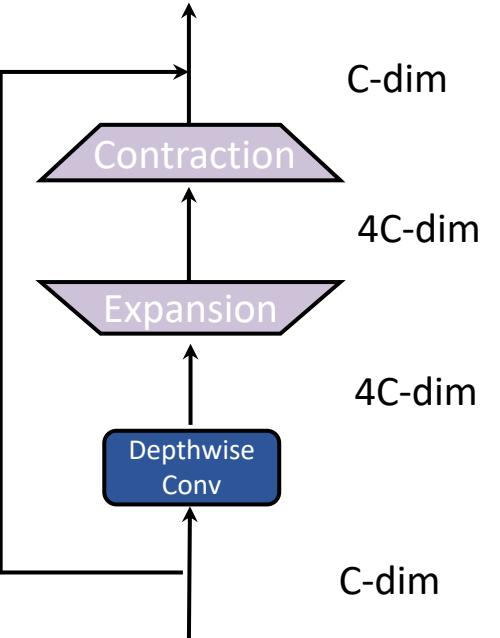
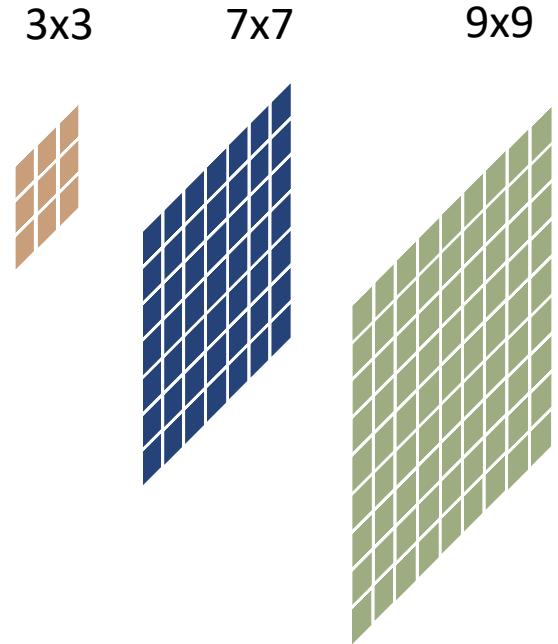
Large kernel size



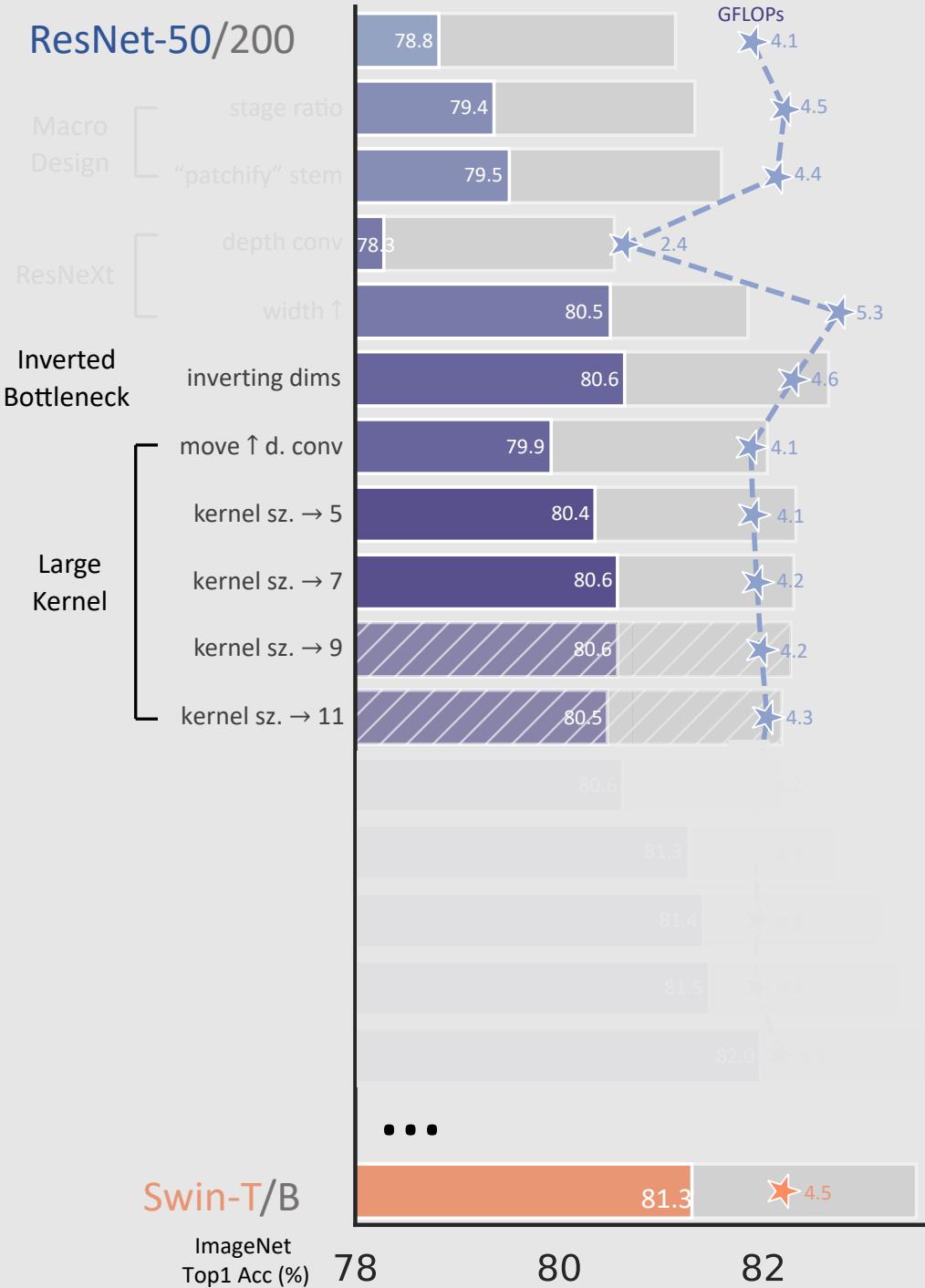
Prerequisite: move depth-wise before “expansion”.
Reduce flops w/ larger kernel size.



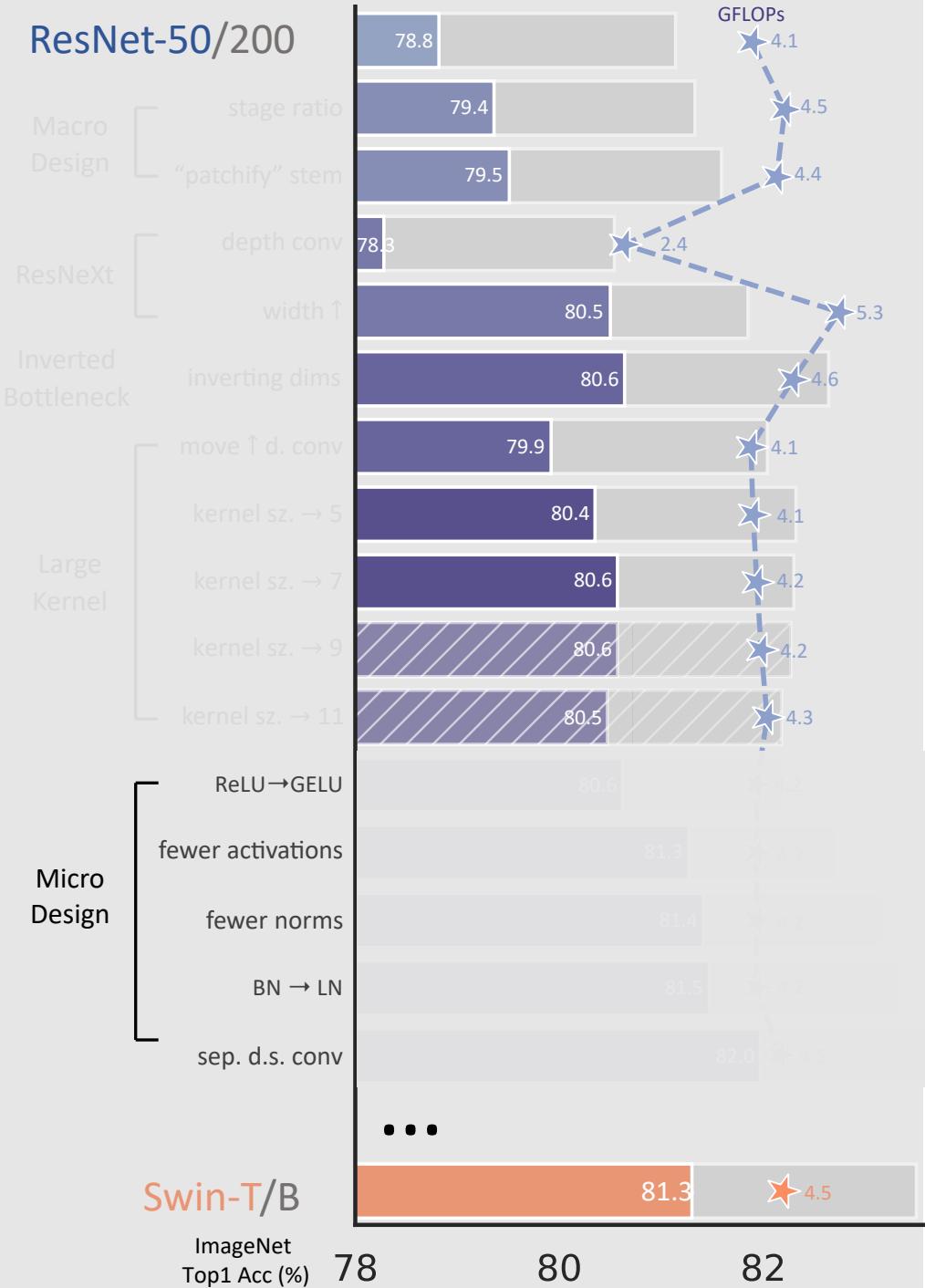
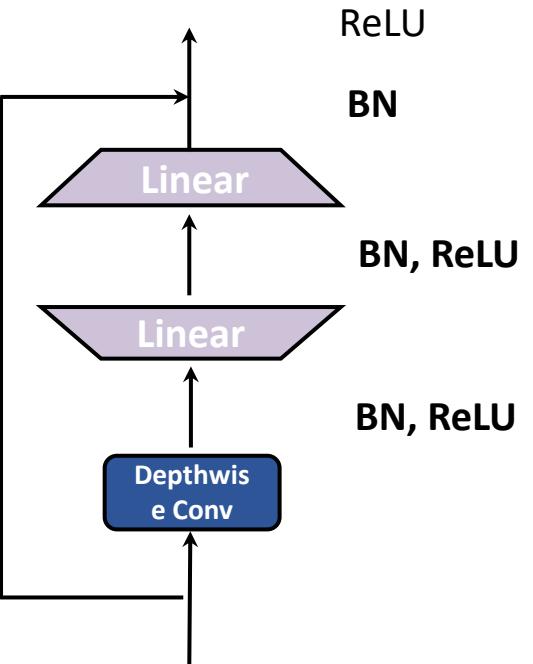
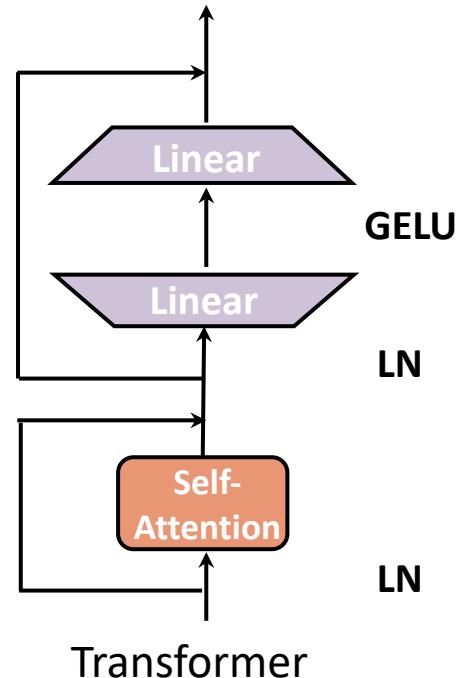
Large kernel size



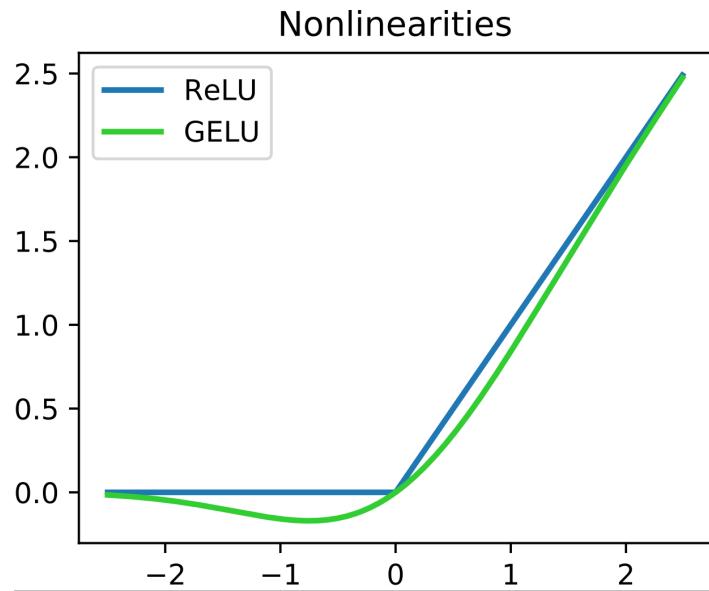
Larger kernel size helps; performance saturates at 7x7
Swin's choice of local window size is also 7 😐



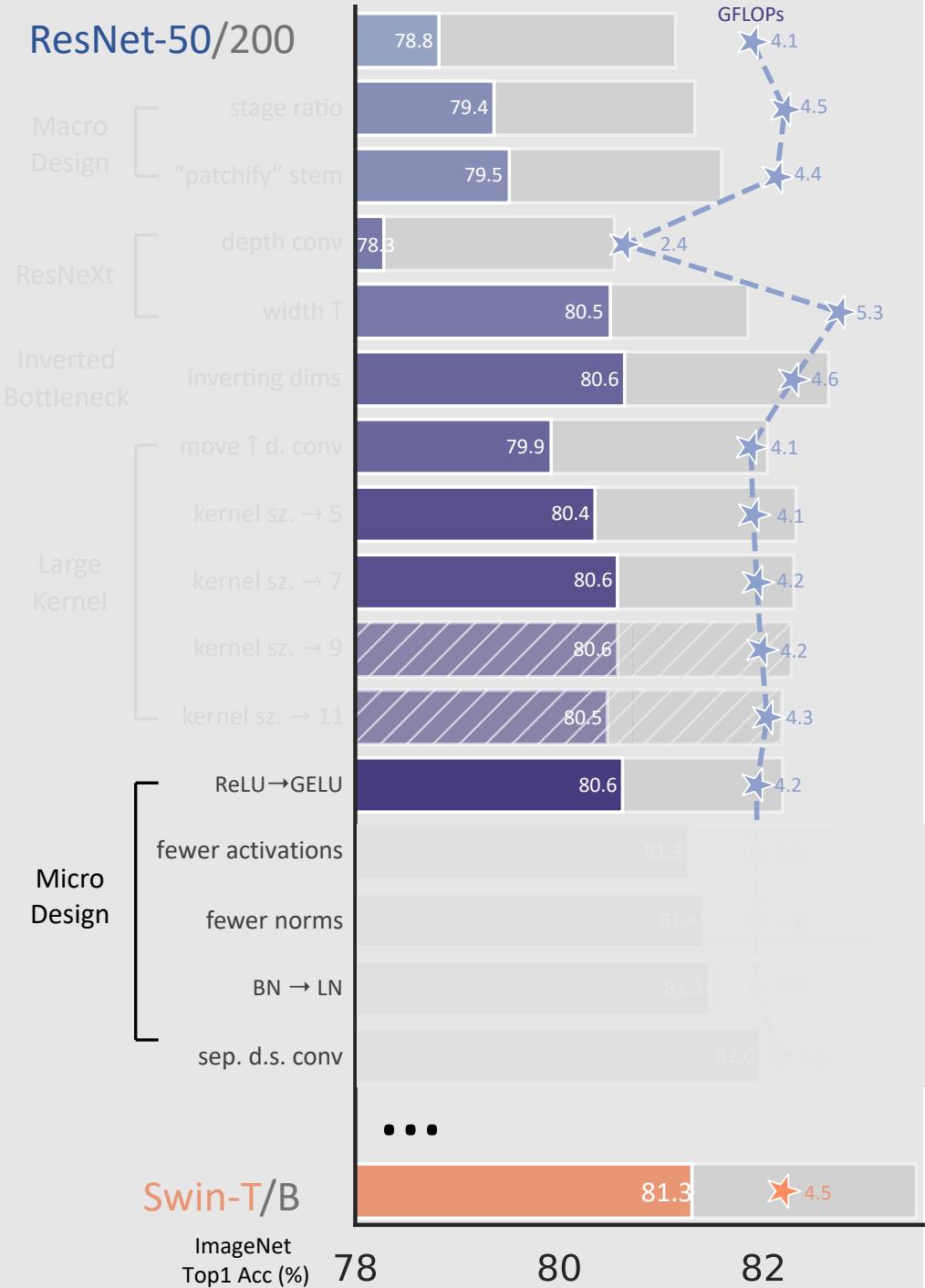
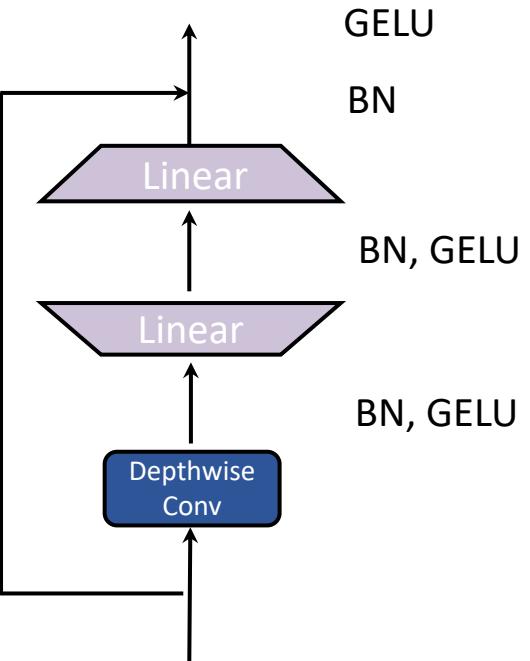
Micro Design



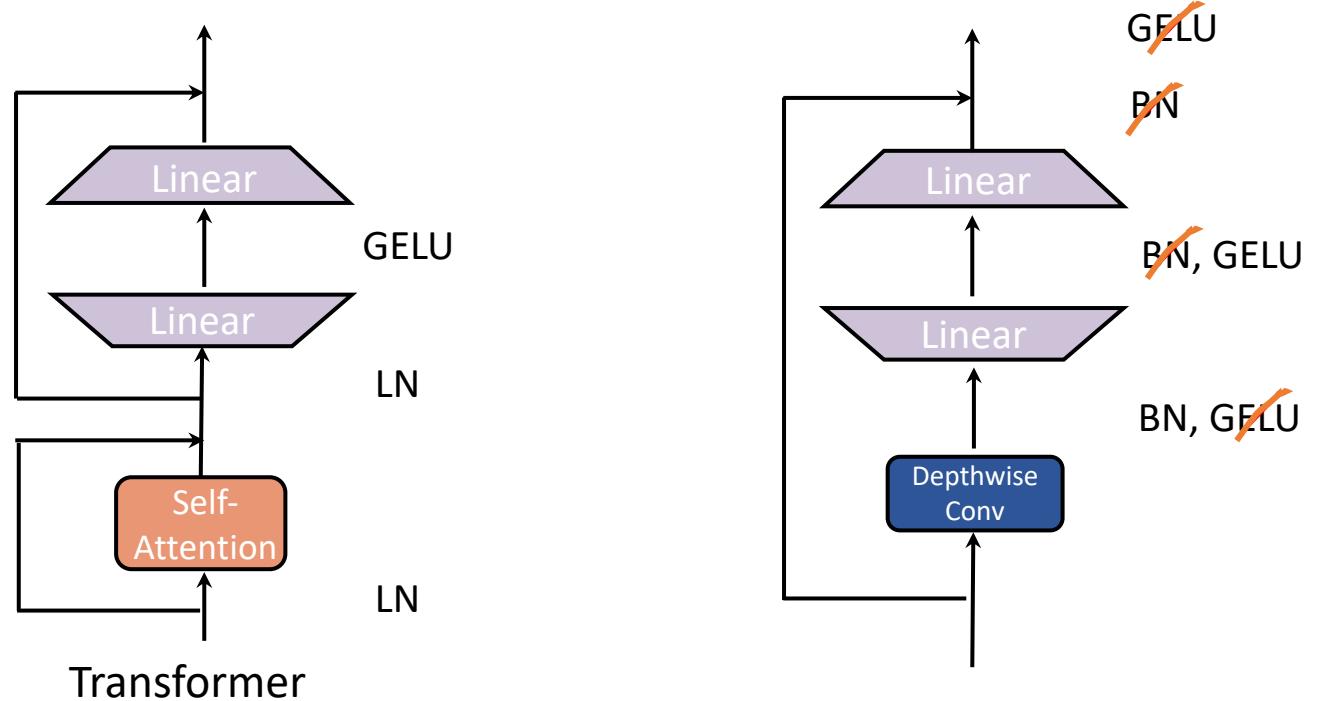
Activations



ReLU → GELU

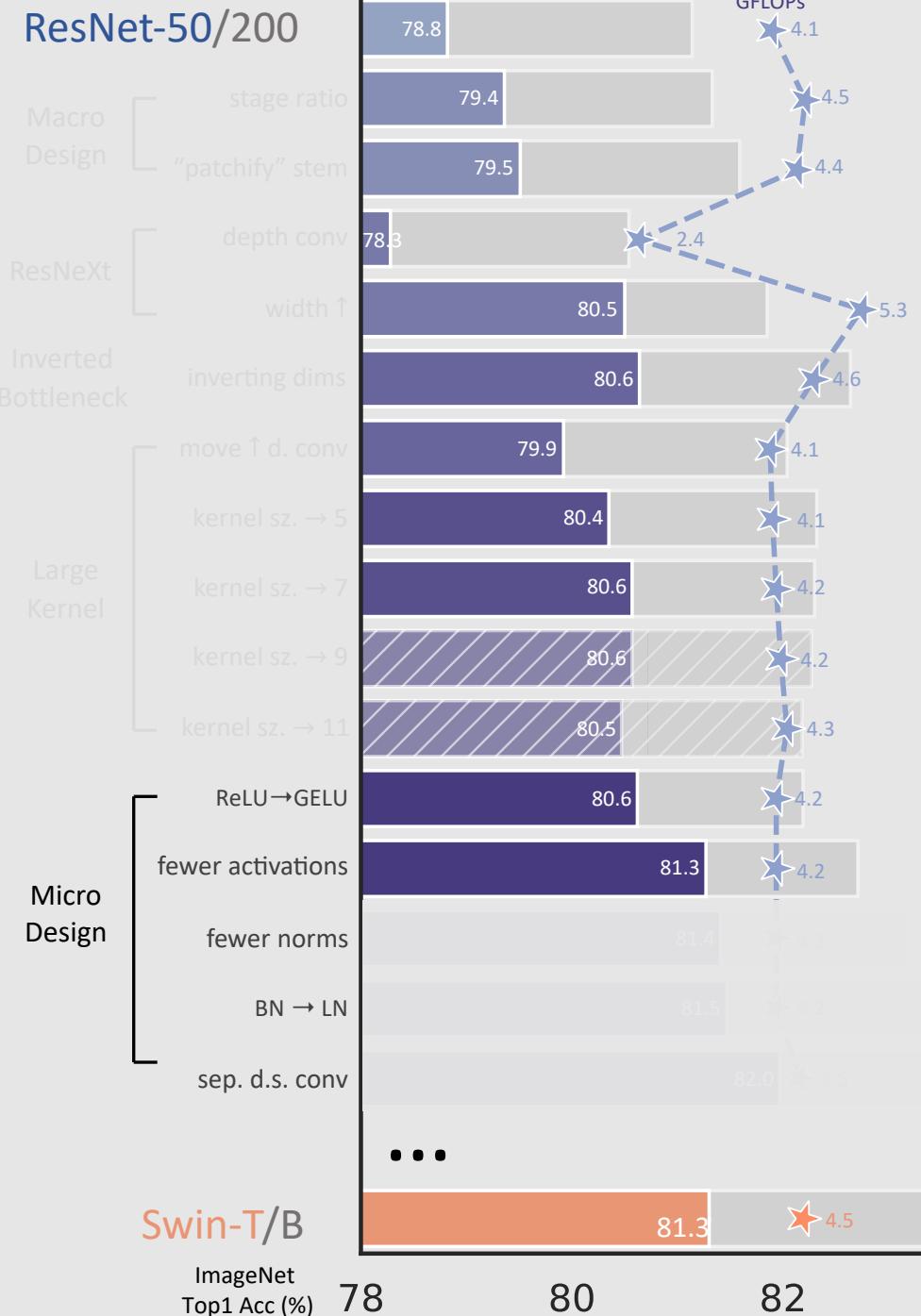


Fewer Activations / Norms



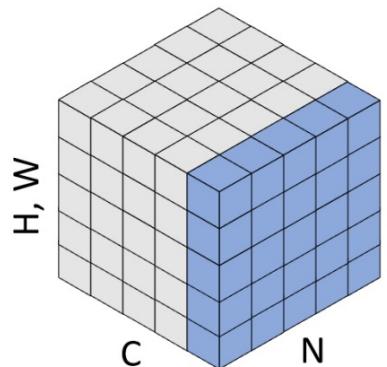
One norm, one activation is good enough per block

+0.8% without changing FLOPs

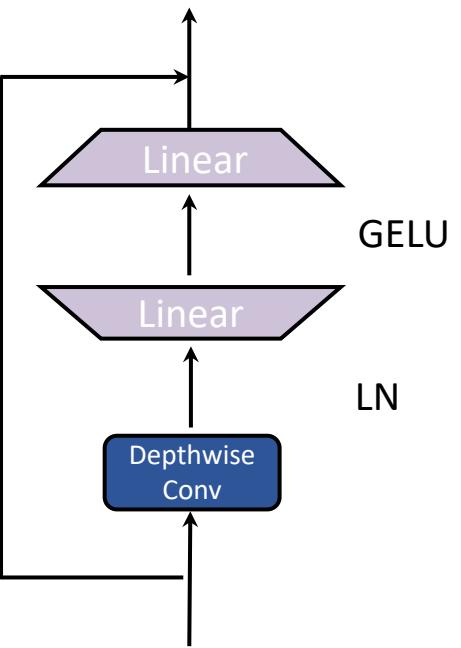
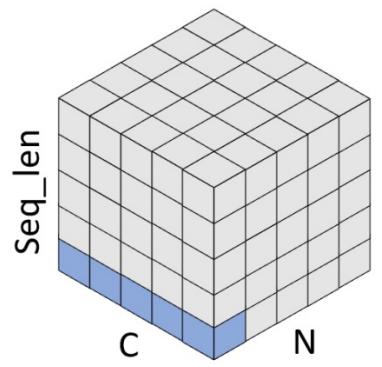


Normalization

BatchNorm



LayerNorm



One **LayerNorm** is good enough.
Say to **BatchNorm** pitfalls!

ResNet-50/200

Macro Design

ResNeXt

Inverted Bottleneck

Large Kernel

Micro Design

Swin-T/B

ImageNet
Top1 Acc (%)

78

80

82

81.3 4.5

81.5 4.2

81.4 4.2

81.3 4.2

80.6 4.2

80.5 4.3

80.6 4.2

80.4 4.1

79.9 4.1

80.6 4.6

80.5 5.3

80.6 4.4

79.5 4.5

78.3 2.4

80.5 4.1

79.4 4.5

78.8 4.1

GFLOPs

4.1

4.5

4.4

2.4

5.3

4.6

4.1

4.1

4.2

4.2

4.3

4.2

4.2

4.2

4.2

4.2

4.2

4.2

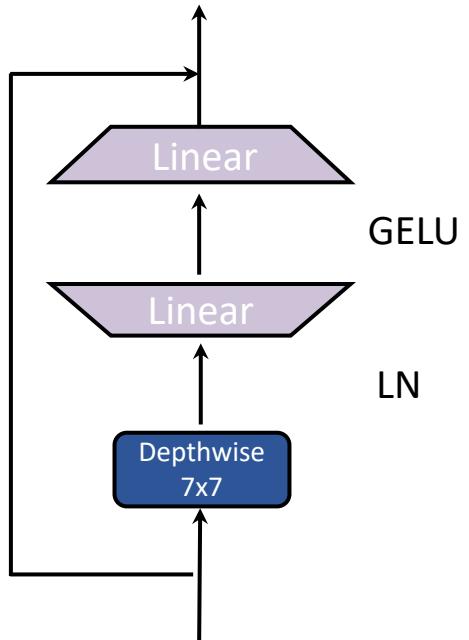
4.2

4.2

4.2

4.5

ConvNeXt Block

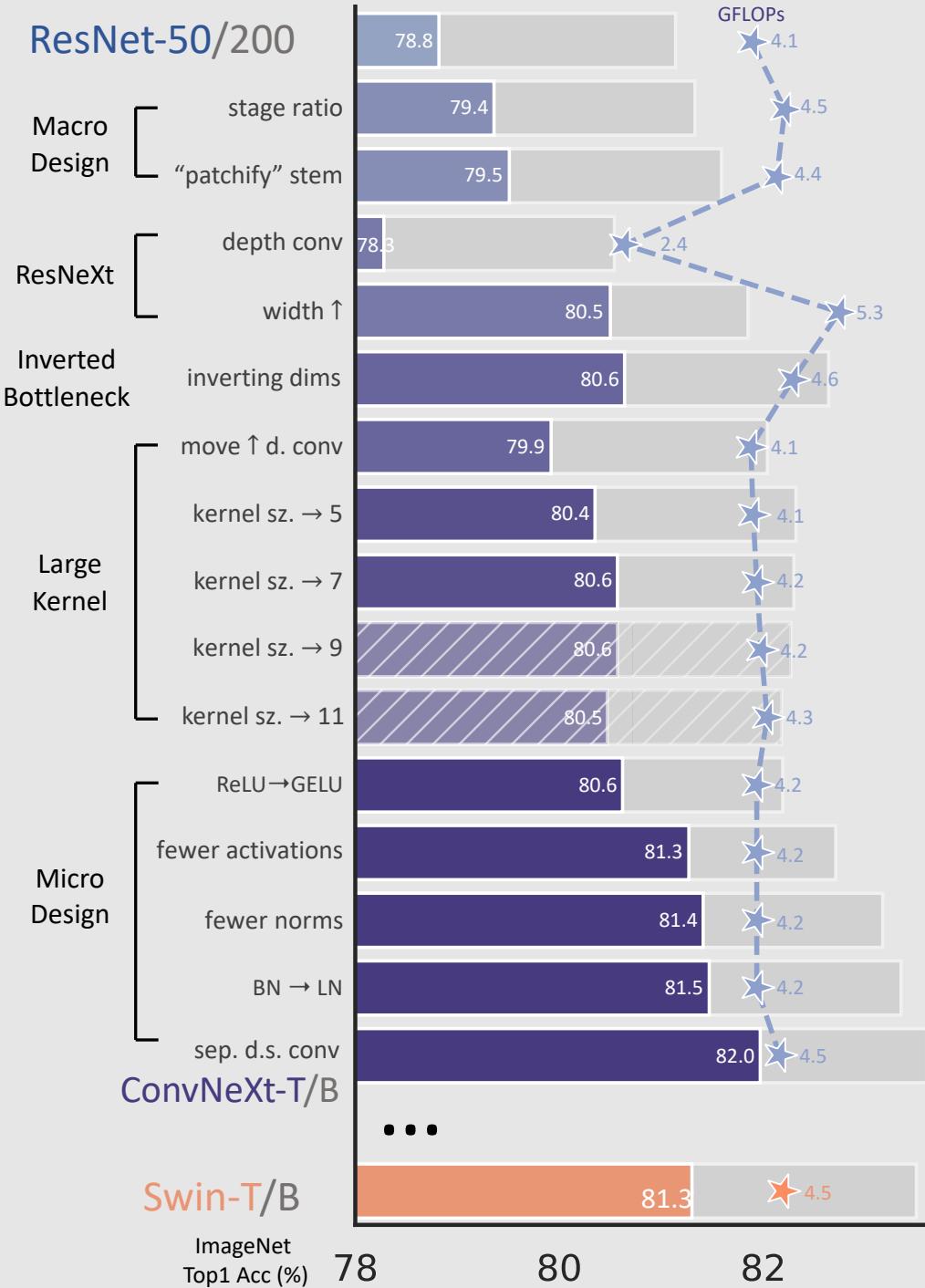


Minimal design:

★ as simple as possible (but not simpler).

★ only 100 lines of code

(vs. 500+ for advanced vision transformers)



```

class Block(nn.Module):
    r"""
    ConvNeXt Block. There are two equivalent implementations:
    (1) DwConv -> LayerNorm (channels_first) -> 1x1 Conv -> GELU -> 1x1 Conv; all in (N, C, H, W)
    (2) DwConv -> Permute to (N, H, W, C); LayerNorm (channels_last) -> Linear -> GELU -> Linear; Permute back
    We use (2) as we find it slightly faster in PyTorch

    Args:
        dim (int): Number of input channels.
        drop_path (float): Stochastic depth rate. Default: 0.0
        layer_scale_init_value (float): Init value for Layer Scale. Default: 1e-6.
    """
    def __init__(self, dim, drop_path=0., layer_scale_init_value=1e-6):
        super().__init__()
        self.dwconv = nn.Conv2d(dim, dim, kernel_size=7, padding=3, groups=dim) # depthwise conv
        self.norm = LayerNorm(dim, eps=1e-6)
        self.pwconv1 = nn.Linear(dim, 4 * dim) # pointwise/1x1 convs, implemented with linear layers
        self.act = nn.GELU()
        self.pwconv2 = nn.Linear(4 * dim, dim)
        self.gamma = nn.Parameter(layer_scale_init_value * torch.ones((dim)),
                                 requires_grad=True) if layer_scale_init_value > 0 else None
        self.drop_path = DropPath(drop_path) if drop_path > 0. else nn.Identity()

    def forward(self, x):
        input = x
        x = self.dwconv(x)
        x = x.permute(0, 2, 3, 1) # (N, C, H, W) -> (N, H, W, C)
        x = self.norm(x)
        x = self.pwconv1(x)
        x = self.act(x)
        x = self.pwconv2(x)
        if self.gamma is not None:
            x = self.gamma * x
        x = x.permute(0, 3, 1, 2) # (N, H, W, C) -> (N, C, H, W)

        x = input + self.drop_path(x)
        return x

```

ConvNeXt block

Swin Transformer block

```

def window_partition(x, window_size):
    """
    Args:
        x: (B, H, W, C)
        window_size (int): window size
    Returns:
        windows: (num_windows*B, window_size, window_size, C)
    """
    B, H, W, C = x.shape
    x = x.view(B, H // window_size, window_size, W // window_size, window_size, C)
    windows = x.permute(0, 1, 3, 2, 4, 5).contiguous().view(-1, window_size, window_size, C)
    return windows

def window_reverse(windows, window_size, H, W):
    """
    Args:
        windows: (num_windows*B, window_size, window_size, C)
        window_size (int): Window size
        H (int): Height of image
        W (int): Width of image
    Returns:
        x: (B, H, W, C)
    """
    B = int(windows.shape[0] / (H * W / window_size / window_size))
    x = windows.view(B, H // window_size, W // window_size, window_size, window_size, C)
    x = x.permute(0, 1, 3, 2, 4, 5).contiguous().view(B, H, W, C)
    return x

class WindowAttention(nn.Module):
    r"""
    Window based multi-head self attention (W-MSA) module with relative position bias. It supports both of shifted and non-shifted window.
    """
    def __init__(self, dim, window_size, num_heads, qkv_bias=False, qk_scale=None,
                 attn_drop=0., proj_drop=0.):
        super().__init__()
        self.dim = dim
        self.window_size = window_size # Wh, Ww
        self.num_heads = num_heads
        head_dim = dim // num_heads
        self.scale = qk_scale or head_dim ** -0.5

        # define a parameter table of relative position bias
        self.relative_position_bias_table = nn.Parameter(
            torch.zeros((2 * window_size[0] - 1) * (2 * window_size[1] - 1), num_heads))

        # get pair-wise relative position index for each token inside the window
        coords_h = torch.arange(self.window_size[0])
        coords_w = torch.arange(self.window_size[1])
        cords = torch.stack(torch.meshgrid([coords_h, coords_w])) # 2, Wh, Ww
        cords_flatten = torch.flatten(cords, 1, 2) # WhWw
        relative_coords = cords_flatten[:, :, None] - cords_flatten[:, None, :]
        relative_coords = relative_coords.permute(0, 2, 1).contiguous() # WhWw, relative_coords[:, :, 0] += self.window_size[0] - 1 # shift to start from 0
        relative_coords[:, :, 0] += self.window_size[1] - 1
        relative_coords[:, :, 1] *= 2 * self.window_size[1] - 1
        relative_position_index = relative_coords.sum(-1) # WhWw, WhWw
        self.register_buffer("relative_position_index", relative_position_index)

        self.qkv = nn.Linear(dim, dim * 3, bias=qkv_bias)
        self.attn_drop = nn.Dropout(attn_drop)
        self.proj = nn.Linear(dim, dim)
        self.proj_drop = nn.Dropout(proj_drop)

        trunc_normal_(self.relative_position_bias_table, std=.02)
        self.softmax = nn.Softmax(dim=-1)

    def forward(self, x, mask=None):
        """
        Args:
            x: input features with shape of (num_windows*B, N, C)
            mask: (0/-inf) mask with shape of (num_windows, Wh*Ww, Wh*Ww) or None
        """
        B, N, C = x.shape
        qkv = self.qkv(x).reshape(B, N, 3, self.num_heads, C // self.num_heads).permute(0, 3, 1, 2, 4)
        q, k, v = qkv[0], qkv[1], qkv[2] # make torchscript happy (cannot use tensor slice on tensor in pytorch script)

        q = q * self.scale
        attn = (q @ k.transpose(-2, -1))

        relative_position_bias = self.relative_position_bias_table[self.relative_position_index.view(-1)].view(
            self.window_size[0] * self.window_size[1], self.window_size[0] * self.window_size[1], -1) # Wh*Ww, Wh*Ww, num_heads
        relative_position_bias = relative_position_bias.permute(2, 0, 1).contiguous() # num_heads, Wh*Ww, Wh*Ww
        attn = attn + relative_position_bias.unsqueeze(0)

        if mask is not None:
            attn = attn + mask.unsqueeze(0)
            attn = attn.view(B, num_heads, N, N) + mask.unsqueeze(1).unsqueeze(2)
            attn = attn.softmax(dim=-1)
            attn = self.softmax(attn)
        else:
            attn = self.softmax(attn)

        attn = self.attn_drop(attn)

        x = (attn @ v).transpose(1, 2).reshape(B, N, C)
        x = self.proj(x)
        x = self.proj_drop(x)
        return x

    def extra_repr(self) -> str:
        return f"dim={self.dim}, window_size={self.window_size}, num_heads={self.num_heads}"

    def flops(self, N):
        # calculate flops for 1 window with token length of N
        flops = 0
        # qkv = self.qkv(x)
        flops += N * self.dim * 3 * self.dim
        # attn = (q @ k.transpose(-2, -1))
        flops += self.num_heads * N * (self.dim // self.num_heads) * N
        # x = (attn @ v)
        flops += self.num_heads * N * N * (self.dim // self.num_heads)
        # x = self.proj(x)
        flops += N * self.dim * self.dim
        return flops

class SwinTransformerBlock(nn.Module):
    r"""
    Swin Transformer Block.
    """
    def __init__(self, dim, input_resolution, num_heads, window_size=7, shift_size=0,
                 img_size=224, patch_size=4, window_ratio=1., qkv_bias=False, qk_scale=None,
                 attn_drop=0., proj_drop=0., act_layer=nn.GELU, norm_layer=nn.LayerNorm):
        super().__init__()
        self.dim = dim
        self.input_resolution = input_resolution
        self.num_heads = num_heads
        self.window_size = window_size
        self.shift_size = shift_size
        self.window_ratio = window_ratio
        self.qkv_bias = qkv_bias
        self.qk_scale = qk_scale
        self.attn_drop = attn_drop
        self.proj_drop = proj_drop
        self.act_layer = act_layer
        self.norm_layer = norm_layer
        self.norm = norm_layer(dim)

        # get mask windows
        mask_windows = window_partition(img_mask, self.window_size) # nhw, window_size, window_size, 1
        mask_windows = mask_windows.view(-1, self.window_size * self.window_size, 1)
        attn_mask = mask_windows.unsqueeze(1) - mask_windows.unsqueeze(2)
        attn_mask = attn_mask.masked_fill(attn_mask == 0, float(-100.0))
        attn_mask = attn_mask + None
        self.register_buffer("attn_mask", attn_mask)

        def forward(self, x):
            H, W = self.input_resolution
            B, L, C = x.shape
            assert L == H * W, "input feature has wrong size"
            shortcut = x
            x = self.norm(x)
            x = x.view(B, H, W, C)

            # cyclic shift
            if self.shift_size > 0:
                shifted_x = torch.roll(x, shifts=(-self.shift_size, -self.shift_size), dims=(1, 2))
            else:
                shifted_x = x

            # partition windows
            X_windows = window_partition(shifted_x, self.window_size) # nhw, window_size, window_size, C
            X_windows = X_windows.view(-1, self.window_size * self.window_size, C) # nhw, window_size*window_size, C
            X_windows = self.norm(X_windows)

            # calculate attention
            if self.window_ratio == 1.:
                attn_windows = self.window_attention(X_windows, self.window_size, num_heads, qkv_bias=qkv_bias, qk_scale=qk_scale, attn_drop=attn_drop, proj_drop=proj_drop)
            else:
                attn_windows = self.window_attention(X_windows, self.window_size, num_heads, qkv_bias=qkv_bias, qk_scale=qk_scale, attn_drop=attn_drop, proj_drop=proj_drop)

            attn_windows = attn_windows.view(-1, self.window_size, self.window_size, C)
            if self.window_ratio == 1.:
                shifted_x = shifted_x + attn_windows
            else:
                shifted_x = shifted_x + self.drop_path(attn_windows)

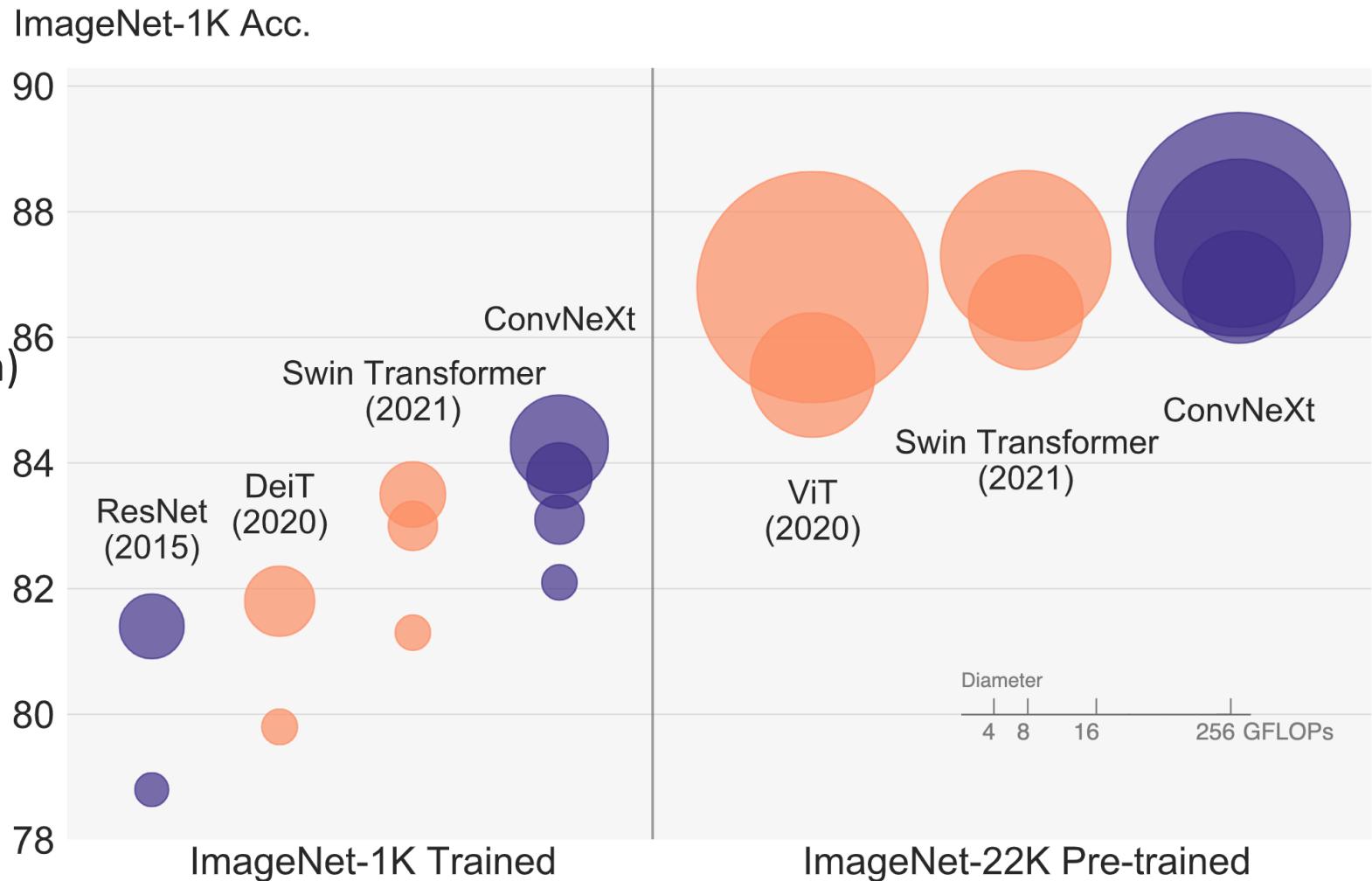
            x = shifted_x
            x = x.view(B, H, W, C)
            x = shortcut + x
            x = x + self.drop_path(self.norm(self.act_layer(x)))

            return x

```

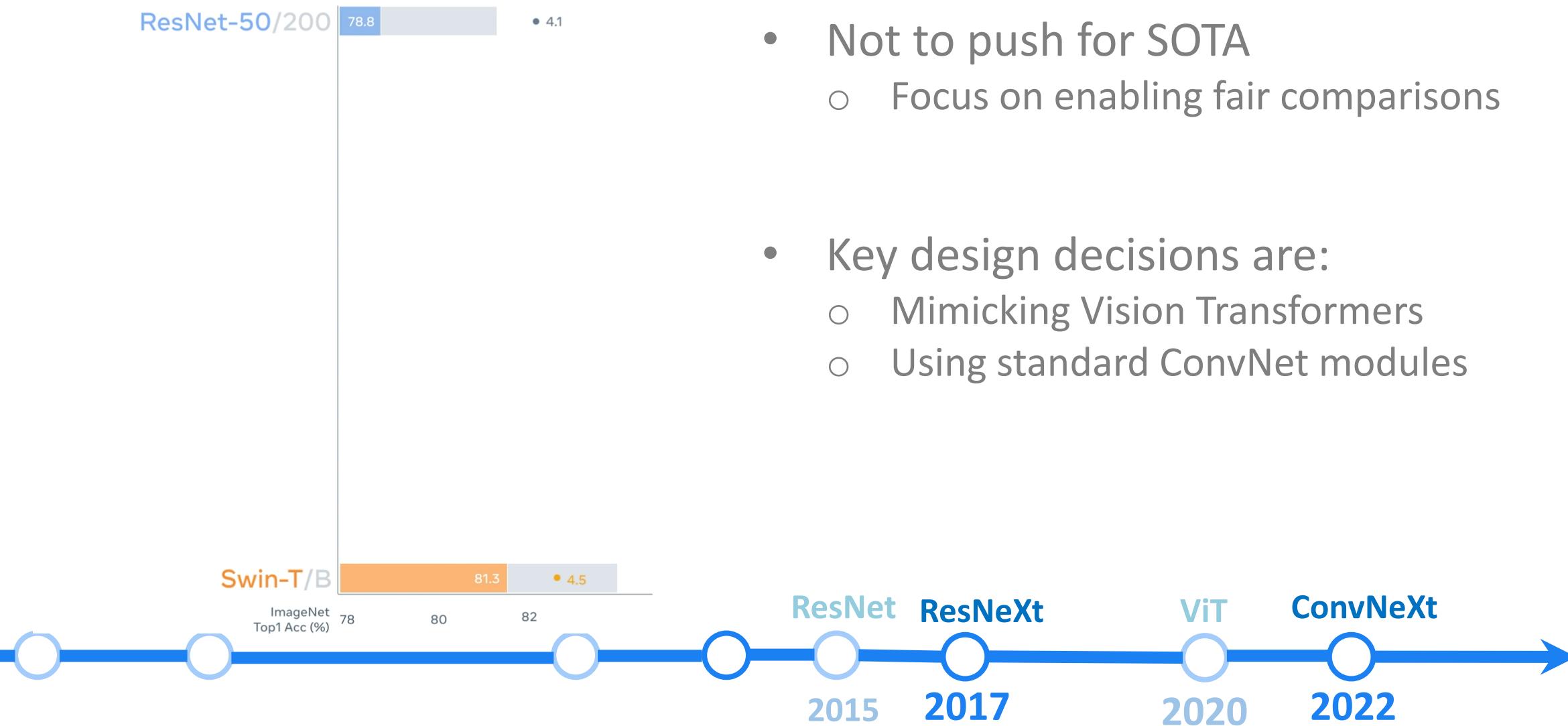
ConvNeXt: Results

- Attention is *NOT* essential
- ConvNets can be **scalable**
(while being much **simpler** in design)



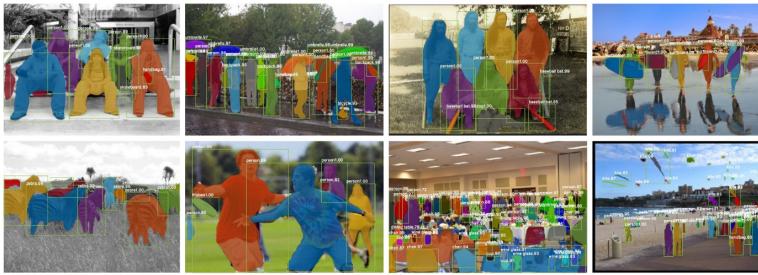
ConvNeXt

[Liu, Mao, Wu, Feichtenhofer, Darrell, Xie. CVPR 2022]



ConvNeXt: Downstream Transfer

Consistently *outperforms* SOTA vision transformers



backbone	FLOPs	FPS	AP _{box}	AP ₅₀	AP ₇₅	AP _{mask}	AP ₅₀	AP ₇₅
Mask-RCNN 3× schedule								
○ Swin-T	267G	23.1	46.0	68.1	50.3	41.6	65.1	44.9
● ConvNeXt-T	262G	25.6	46.2	67.9	50.8	41.7	65.0	44.9
Cascade Mask-RCNN 3× schedule								
● ResNet-50	739G	11.4	46.3	64.3	50.5	40.1	61.7	43.4
● X101-32	819G	9.2	48.1	66.5	52.4	41.6	63.9	45.2
● X101-64	972G	7.1	48.3	66.4	52.3	41.7	64.0	45.1
○ Swin-T	745G	12.2	50.4	69.2	54.7	43.7	66.6	47.3
● ConvNeXt-T	741G	13.5	50.4	69.1	54.8	43.7	66.5	47.3
○ Swin-S	838G	11.4	51.9	70.7	56.3	45.0	68.2	48.8
● ConvNeXt-S	827G	12.0	51.9	70.8	56.5	45.0	68.4	49.1
○ Swin-B	982G	10.7	51.9	70.5	56.4	45.0	68.1	48.9
● ConvNeXt-B	964G	11.4	52.7	71.3	57.2	45.6	68.9	49.5
○ Swin-B [‡]	982G	10.7	53.0	71.8	57.5	45.8	69.4	49.7
● ConvNeXt-B [‡]	964G	11.5	54.0	73.1	58.8	46.9	70.6	51.3
○ Swin-L [‡]	1382G	9.2	53.9	72.4	58.8	46.7	70.1	50.8
● ConvNeXt-L [‡]	1354G	10.0	54.8	73.8	59.8	47.6	71.3	51.7
● ConvNeXt-XL [‡]	1898G	8.6	55.2	74.2	59.9	47.7	71.6	52.2

COCO Detection and Instance Segmentation



backbone	input crop.	mIoU	#param.	FLOPs
ImageNet-1K pre-trained				
○ Swin-T	512 ²	45.8	60M	945G
● ConvNeXt-T	512 ²	46.7	60M	939G
○ Swin-S	512 ²	49.5	81M	1038G
● ConvNeXt-S	512 ²	49.6	82M	1027G
○ Swin-B	512 ²	49.7	121M	1188G
● ConvNeXt-B	512 ²	49.9	122M	1170G
ImageNet-22K pre-trained				
○ Swin-B [‡]	640 ²	51.7	121M	1841G
● ConvNeXt-B [‡]	640 ²	53.1	122M	1828G
○ Swin-L [‡]	640 ²	53.5	234M	2468G
● ConvNeXt-L [‡]	640 ²	53.7	235M	2458G
● ConvNeXt-XL [‡]	640 ²	54.0	391M	3335G

ADE20K Semantic Segmentation

ConvNeXt

[Liu, Mao, Wu, Feichtenhofer, Darrell, Xie. CVPR 2022]

ResNet-50/200 78.8 • 4.1

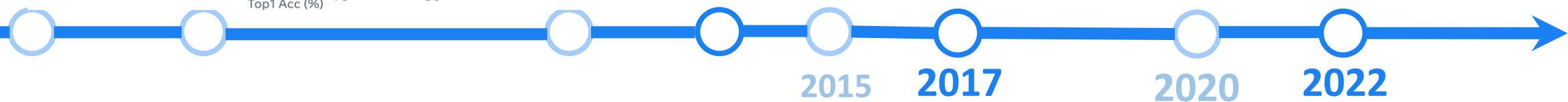
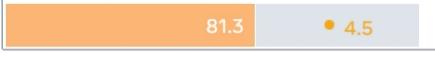
- Not to push for SOTA
 - Focus on enabling fair comparisons



ConvNeXt represents the **community effort!**

Many design choices have been examined **separately** over the last decade, but *not collectively*.

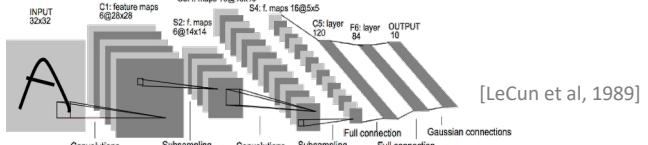
Swin-T/B
ImageNet
Top1 Acc (%)



Connectivity pattern matters

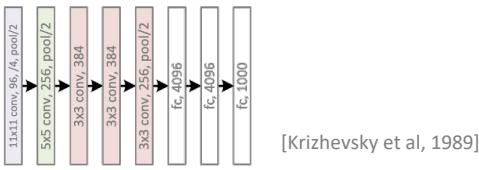
chain-like wiring patterns

LeNet



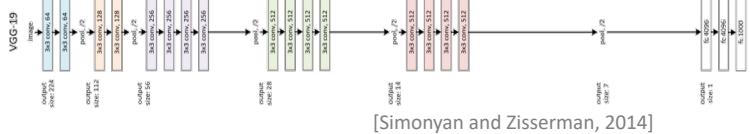
[LeCun et al., 1989]

AlexNet



[Krizhevsky et al., 1989]

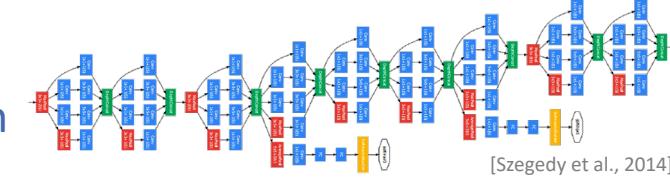
VGGNet



[Simonyan and Zisserman, 2014]

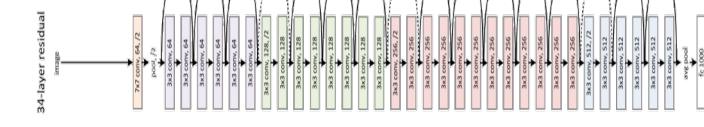
connectivity pattern = graph structure

Inception



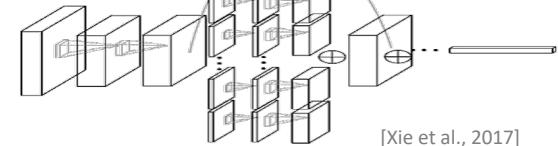
[Szegedy et al., 2014]

ResNet



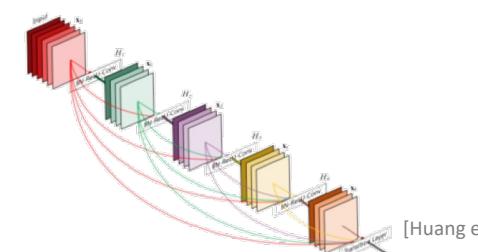
[He et al., 2015]

ResNeXt



[Xie et al., 2017]

DenseNet

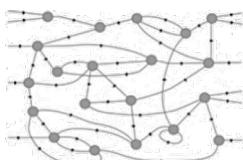
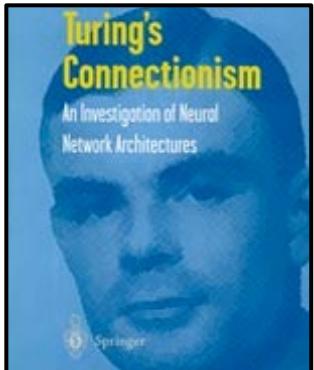


[Huang et al., 2017]

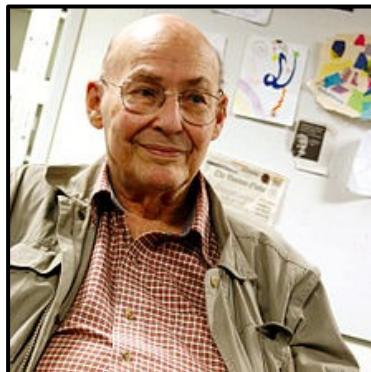
multiple wiring paths

How to wire your computational networks?

Turing (1948)
“Unorganized machines”

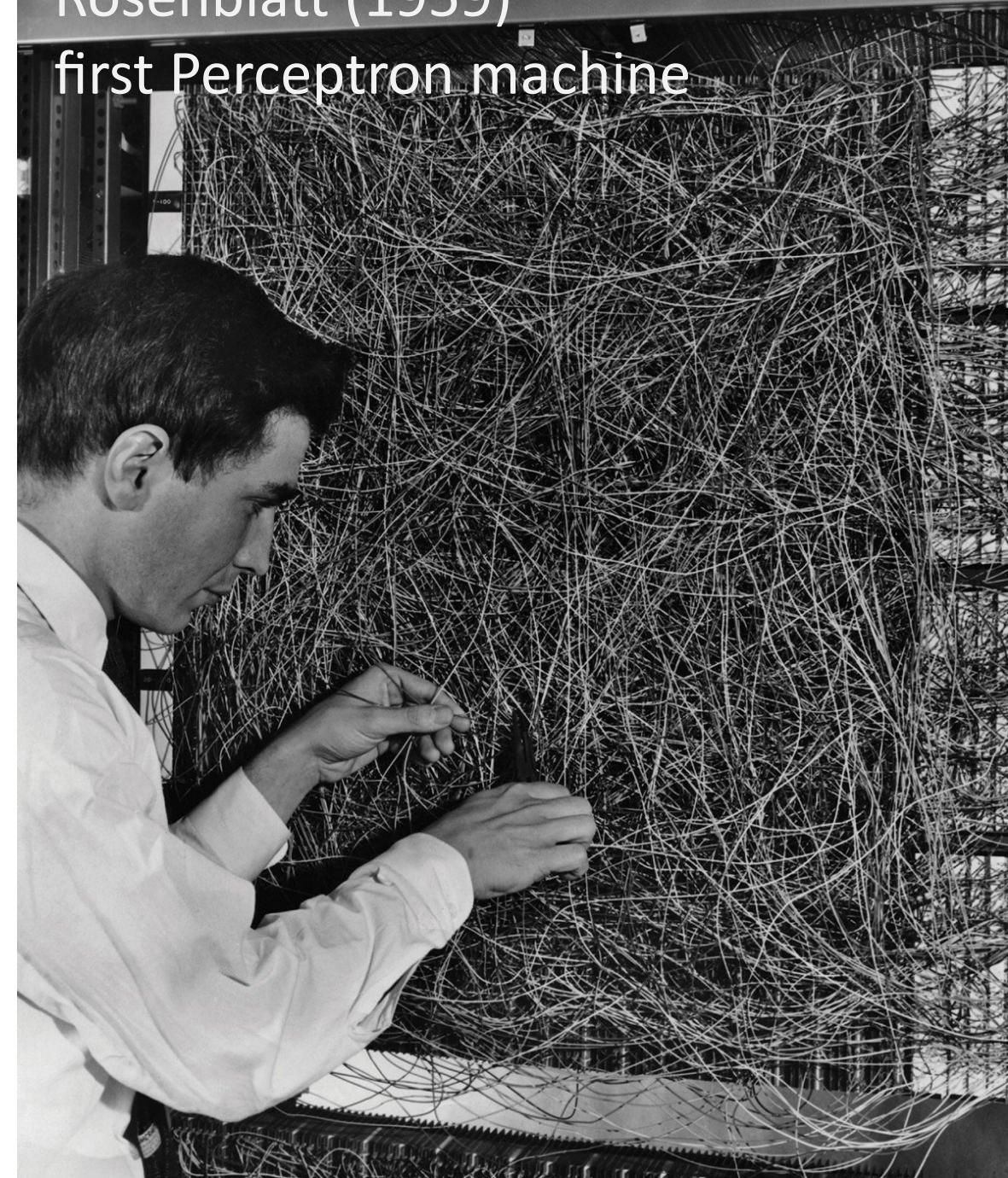


Minsky (1951)
randomly wired “SNARC”

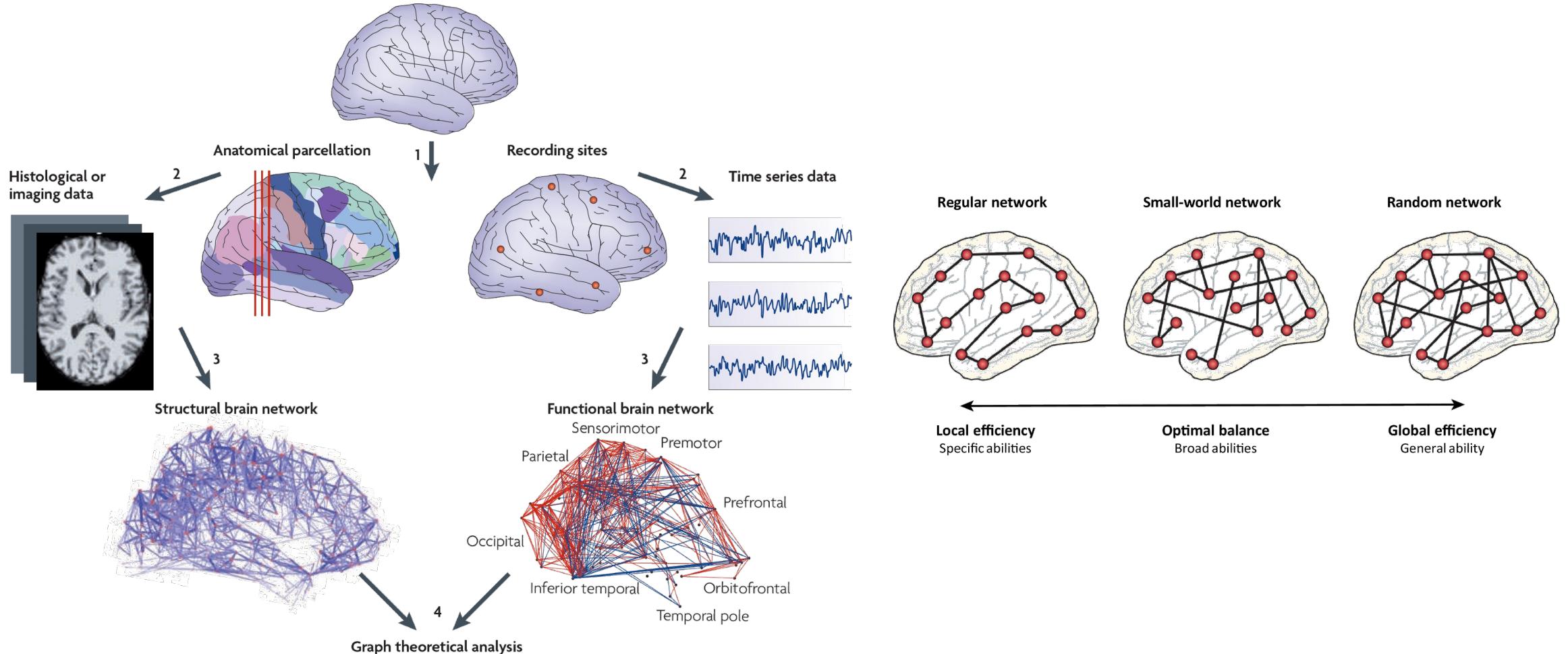


Rosenblatt (1959)

first Perceptron machine



Neuroscience: brain networks are “small-world”?

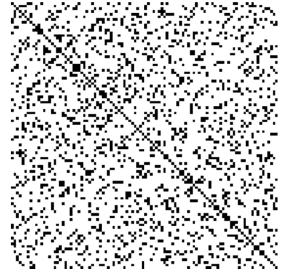
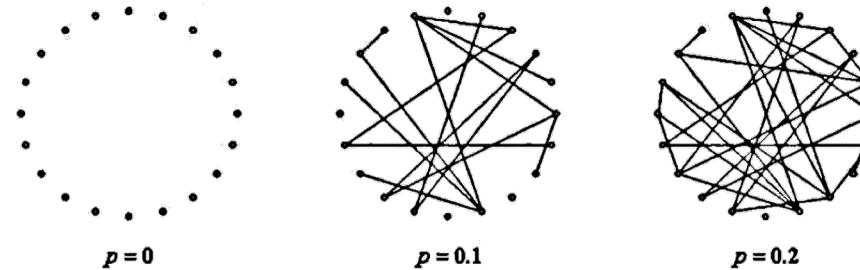


[Bassett, D.S. and Bullmore, E.D., 2006. Small-world brain networks. *The neuroscientist*]

[Bullmore, E. and Sporns, O., 2012. The economy of brain network organization. *Nature reviews neuroscience*]

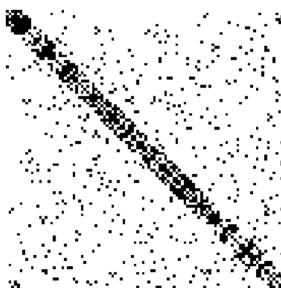
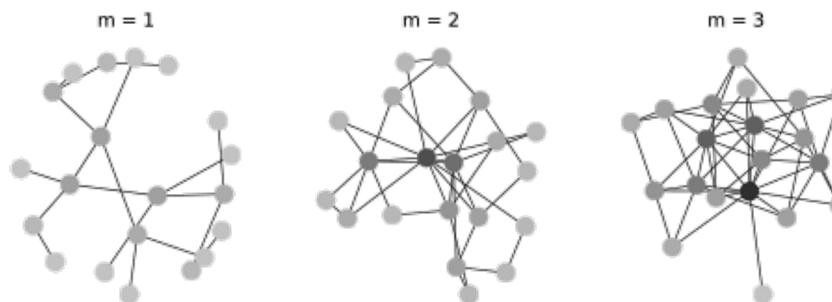
Graph Theory

Erdös-Rényi (ER), 1959



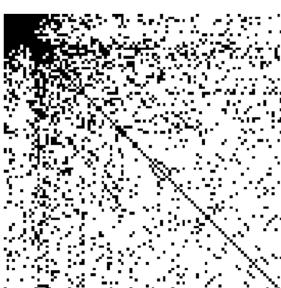
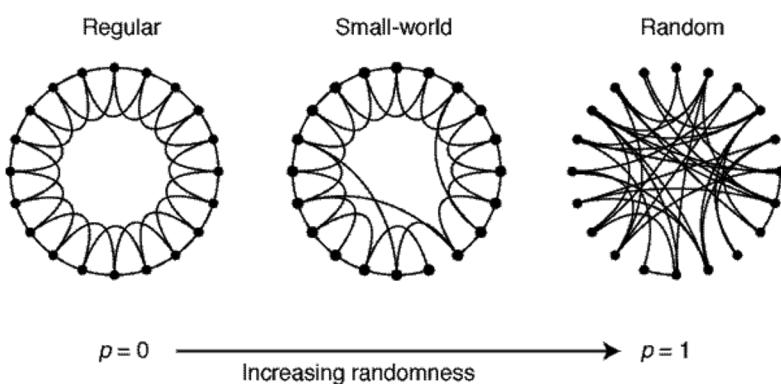
Barabási–Albert (BA), 1998

“Scale-free” graphs



Watts–Strogatz (WS), 1998

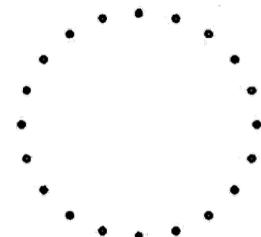
“Small-world” graphs



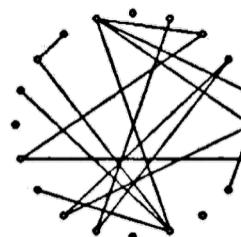
Random Graph Models

Erdös-Rényi (ER), 1959

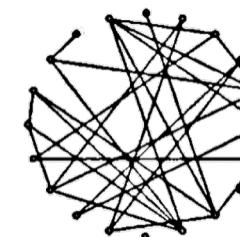
- $G(n, p)$ model: n nodes; each edge is included in the graph with probability p
- Degree distribution: $P(\deg(v) = k) = \binom{n-1}{k} p^k (1-p)^{n-1-k}$
- Almost surely a single connected component if $p > \frac{\ln(n)}{n}$



$p = 0$



$p = 0.1$

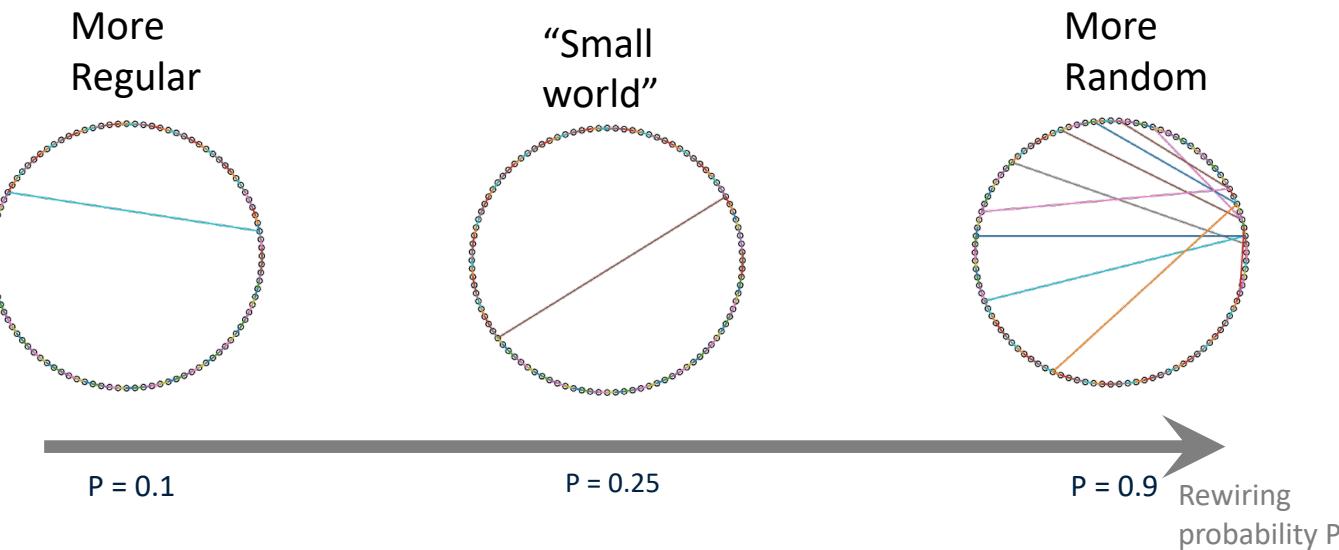


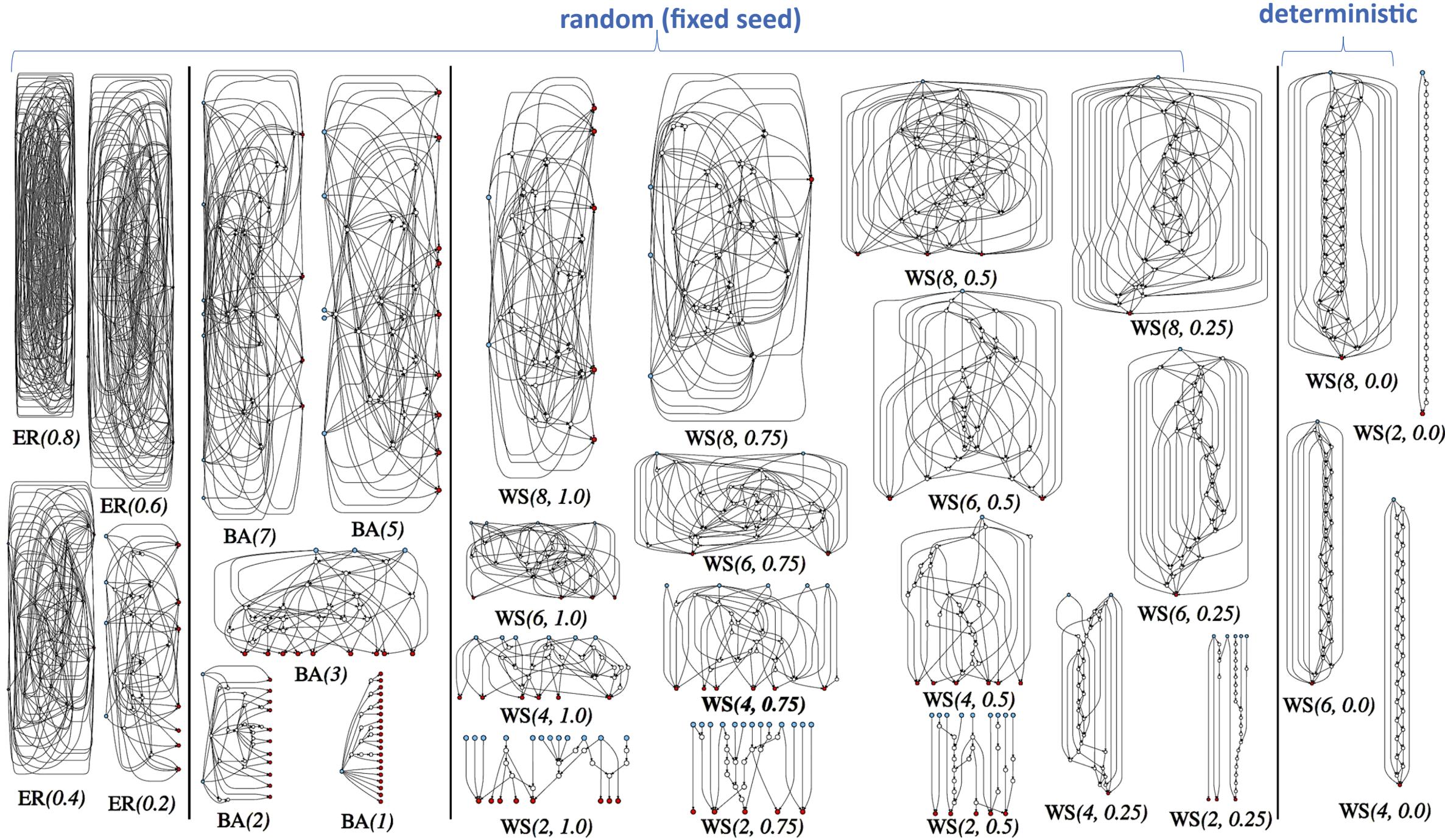
$p = 0.2$

Random Graph Models

Watts–Strogatz (WS), 1998

- “*Small World*” model: high clustering; formation of “Hubs”.
- N nodes regularly placed in a ring; connected to its $K/2$ neighbors on both sides.
- Then each edge is **rewired** with probability p . Randomness enables “*shortcuts*”.





random (fixed seed)

deterministic

RandWire

[Xie, Kirillov, Girshick, He. ICCV 2019]

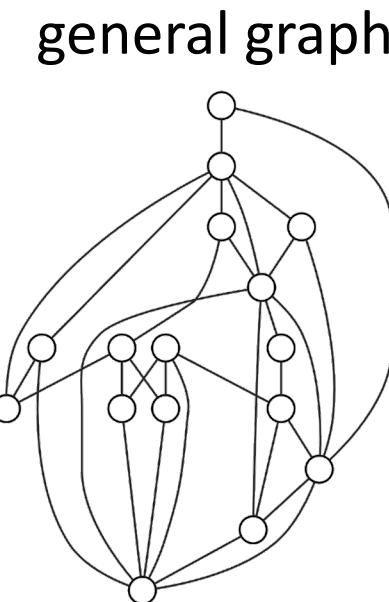
Can we use these graphs to generate trainable neural networks?

Step 1. Generating a general graph G .

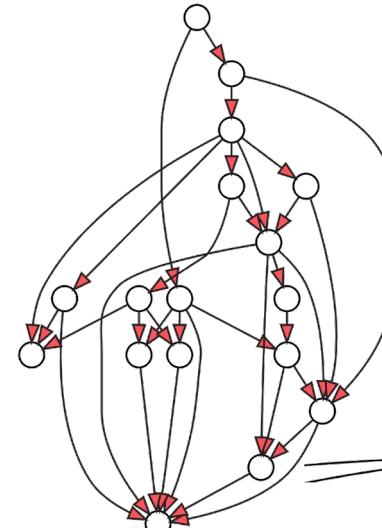
Step 2. Map graph G to a *computable* neural network.

RandWire

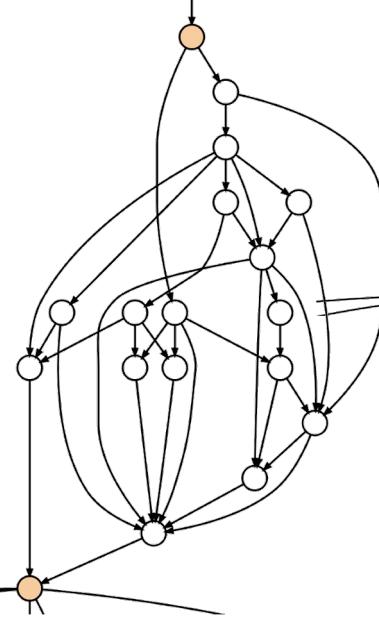
[Xie, Kirillov, Girshick, He. ICCV 2019]



convert to DAG

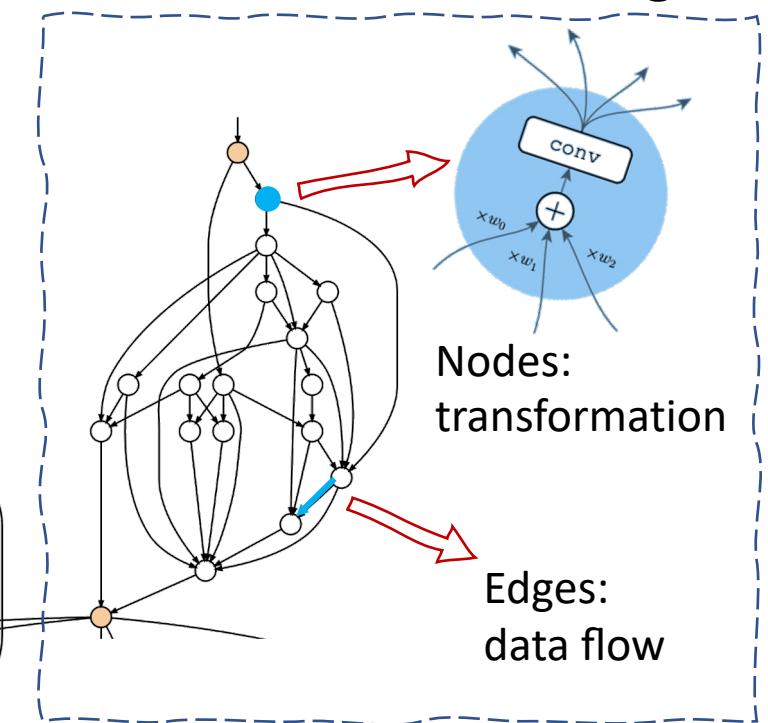


add extra input & output node

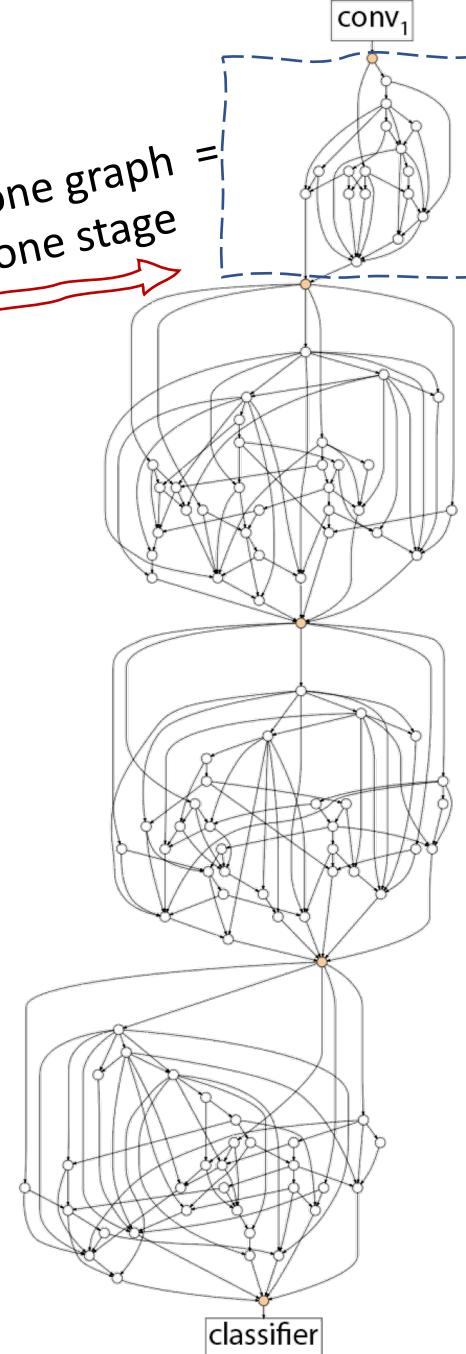


From a graph to a neural network

what are nodes and edges?



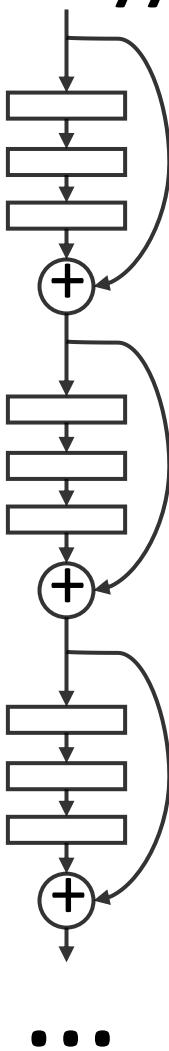
one graph =
one stage



ResNet-50

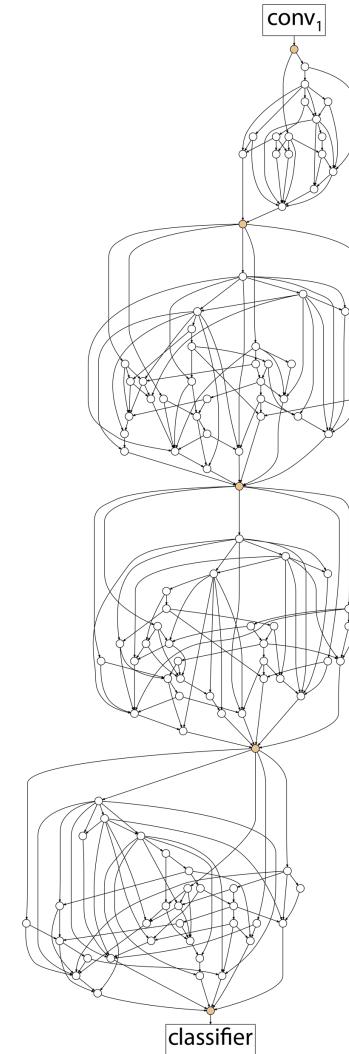
ImageNet Top1 Acc:

77.1%

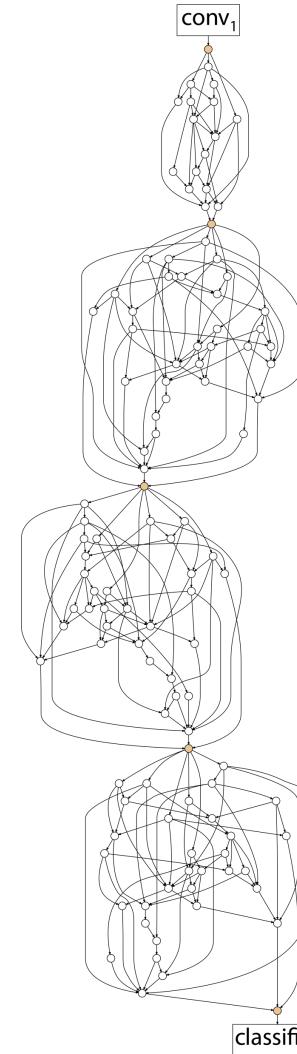


RandWire (same FLOPs)

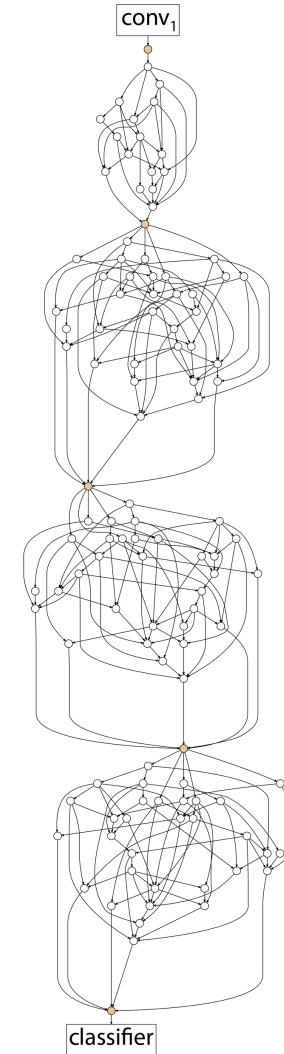
79.1%



79.1%



79.0%



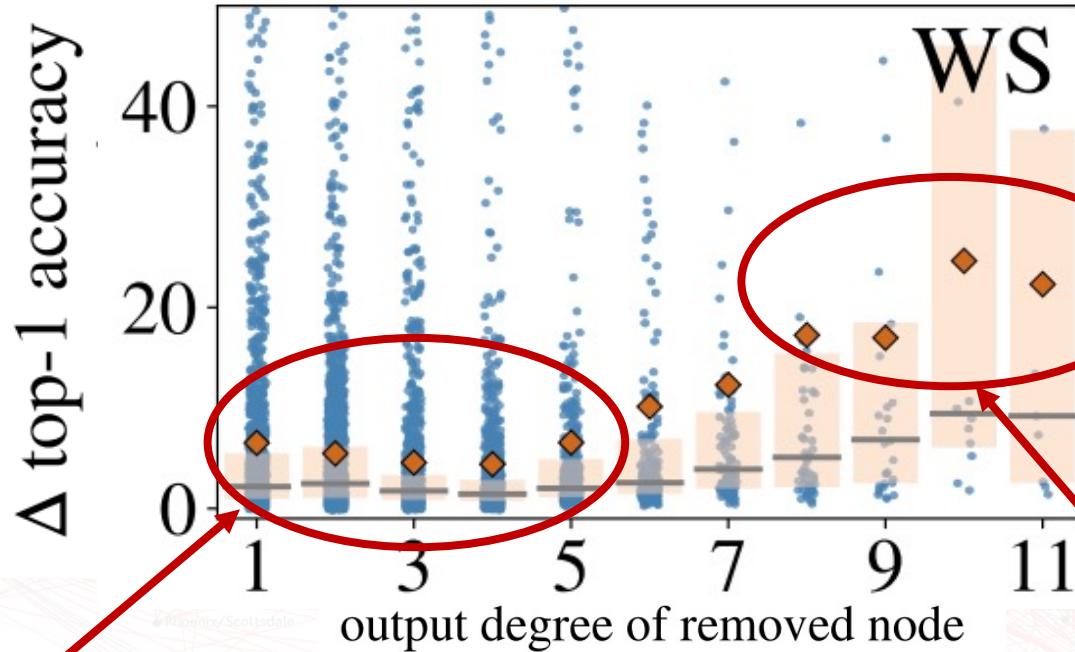
Results

Network	Top-1 Acc	FLOPs
MobileNet	70.6	569
MobileNet v2	74.7	585
ShuffleNet	73.7	524
ShuffleNet v2	74.9	591
NASNet-A	74.0	564
NASNet-C	72.5	558
AmoebaNet-A	74.5	555
AmoebaNet-C	75.7	570
PNAS	74.2	588
DARTS	73.1	595
RandWire-WS	74.7_{±0.25}	583_{±6.2}

Robustness of small-world graph



Graph Damage experiment



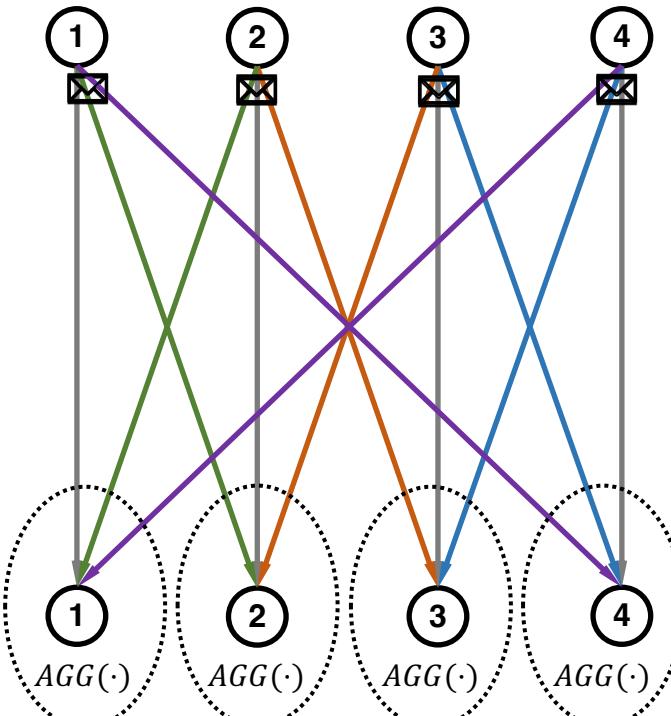
Delete a *non-hub* node:
 < 10% accuracy drop

Delete a “*hub*” node:
 20%-30% accuracy drop

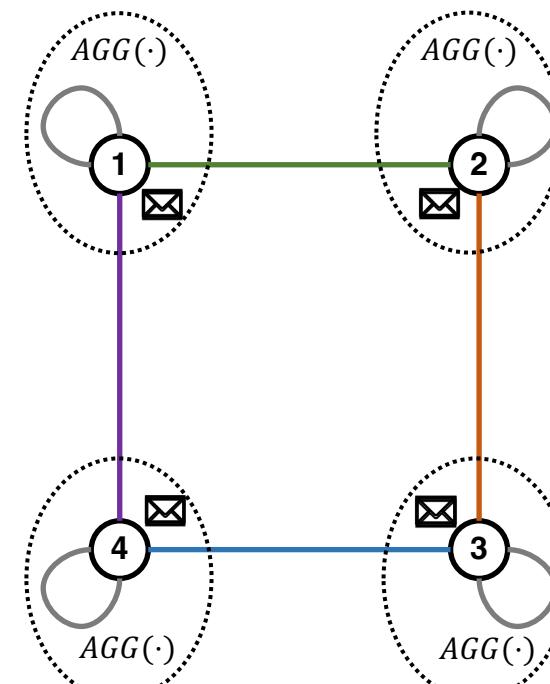
Neuron-level wiring: relational graph

[You, Leskovec, He and Xie. ICML 2020]

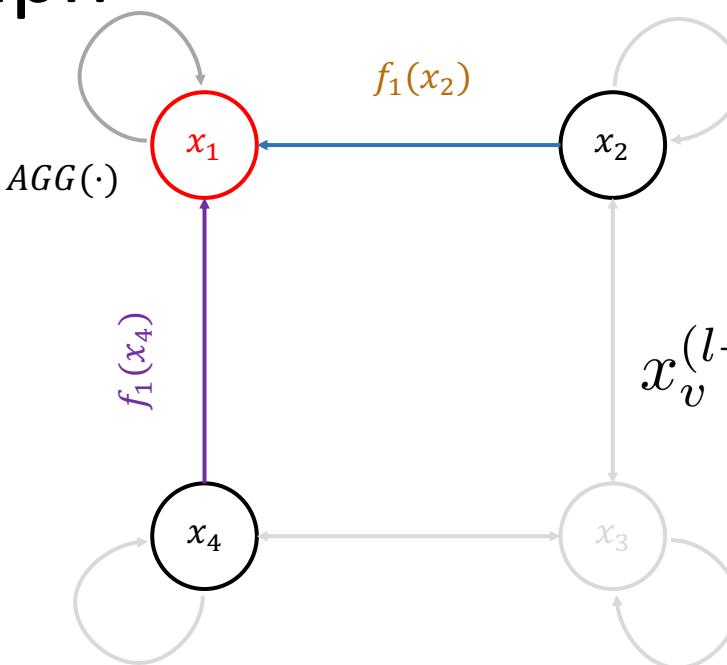
Computational graphs *Directed message flow*



Relational graphs *Bi-Directional message exchange*



Relational graph



$$x_v^{(l+1)} = AGG^{(l)}(\{f_v^{(l)}(x_u^{(l)}), \forall u \in N(v)\})$$

4 Key components	Fixed-width MLP	Variable-width MLP	ResNet-34
Node feature \mathbf{x}_i	Scalar: 1 dimension of data	Vector: multiple dimensions of data	Tensor: multiple channels of data
Message function $f_i(\cdot)$	Scalar multiplication	(Non-square) matrix multiplication	3×3 Conv
Aggregation function $AGG(\cdot)$	$\sigma(\sum(\cdot))$	$\sigma(\sum(\cdot))$	$\sigma(\sum(\cdot))$
Number of rounds L	1 round per layer	1 round per layer	34 rounds with residual connections

Graph structure measures

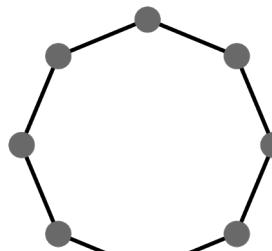
average path length (L)

$$L = \sum_{s,t \in V} \frac{\text{shortest_path}(s, t)}{n(n - 1)}$$

clustering coefficient (C)

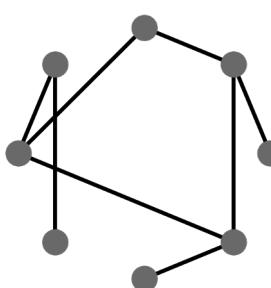
$$C_u = \frac{2T(u)}{\deg(u)(\deg(u) - 1)}$$

WS graph
($n=8, k=2, p=0$)



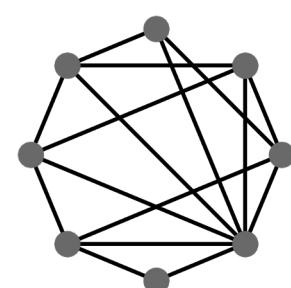
$$L = 2.3 \\ C = 0$$

WS graph
($n=8, k=2, p=0.5$)



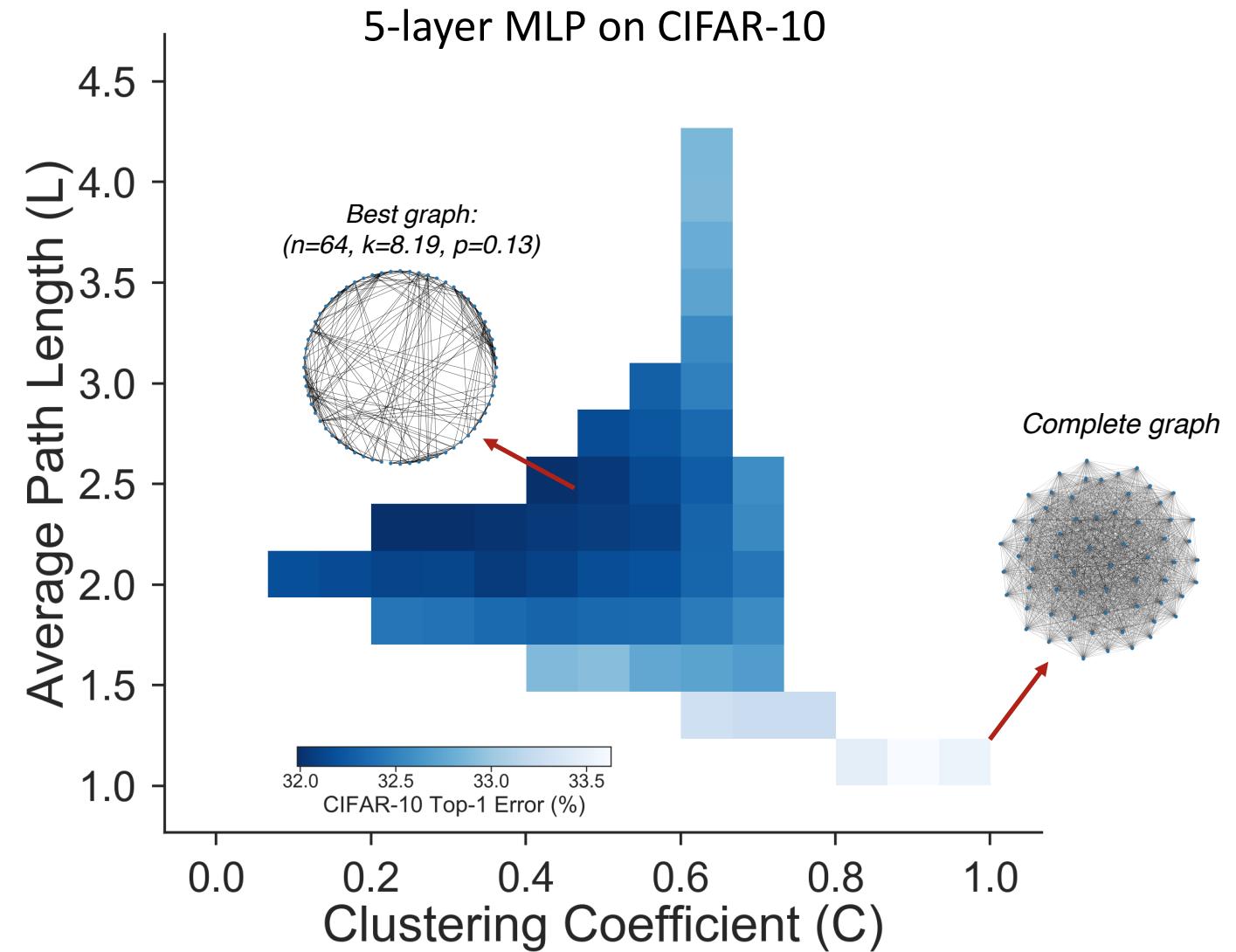
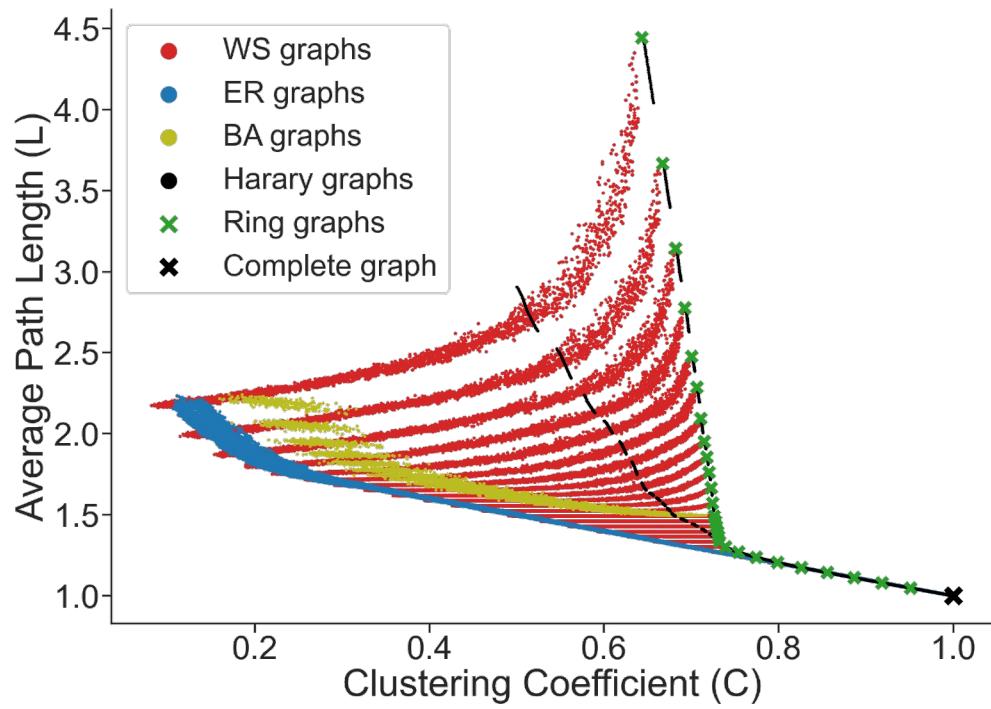
$$L = 2.8 \\ C = 0.3$$

WS graph
($n=8, k=4, p=0.5$)



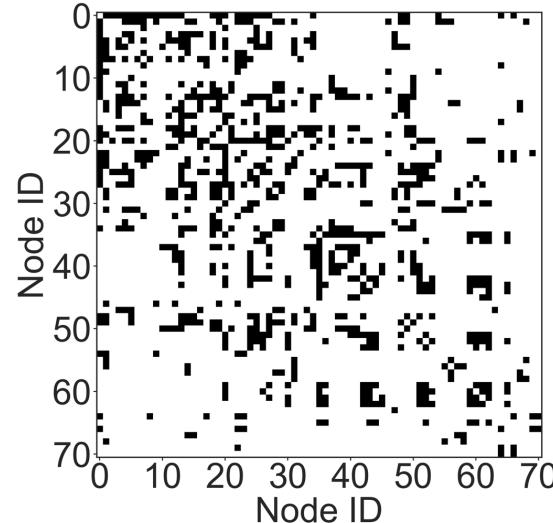
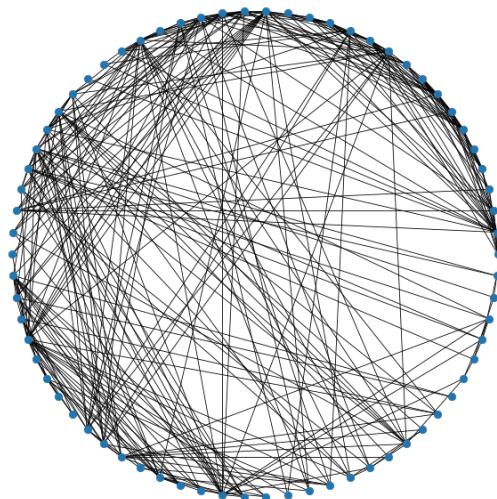
$$L = 1.4 \\ C = 0.4$$

Graph structure measures



Biological neural network: *Macaque whole cortex*

Data from anatomical connectivity matrix [Bassett & Bullmore, 2006]

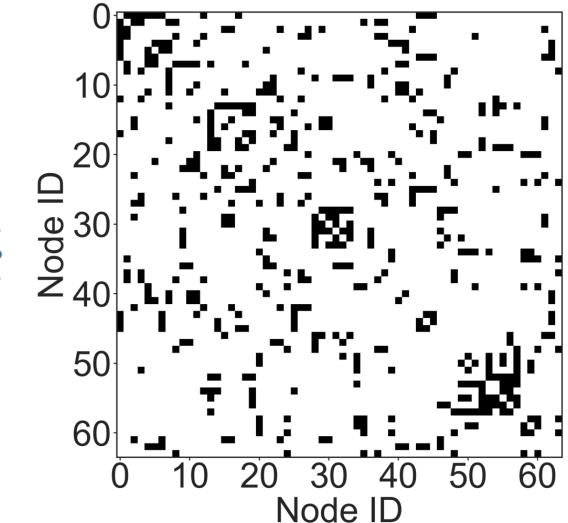


$$L = 2.38$$

$$C = 0.46$$

CIFAR-10: +0.5%

Artificial neural network: *Best 5-layer MLP*



$$L = 2.48$$

$$C = 0.45$$

CIFAR-10: +1.2%

Baseline complete graph:
$L = 1.0$
$C = 1.0$
CIFAR-10: 66.7%

Graph structure of NN

- Graph theory is a *common language*.
- Interests from neuroscience, physics, network science, complex systems, ...



[Analyzing Neural Networks Based on Random Graphs, Janik, and Nowak, arXiv 2020]

[Van Essen et al., 2013]



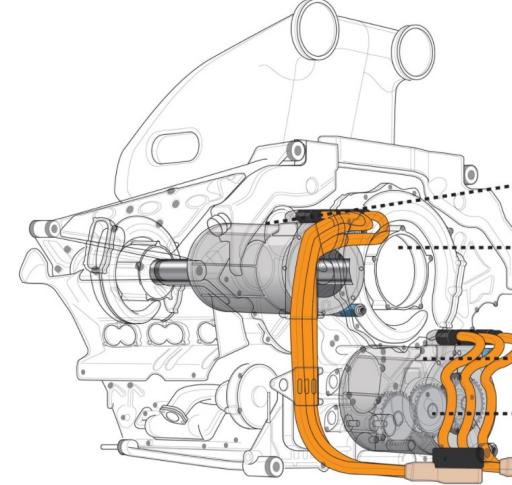
Takeaways: 1. Neural network design is important
2. Wiring pattern matters



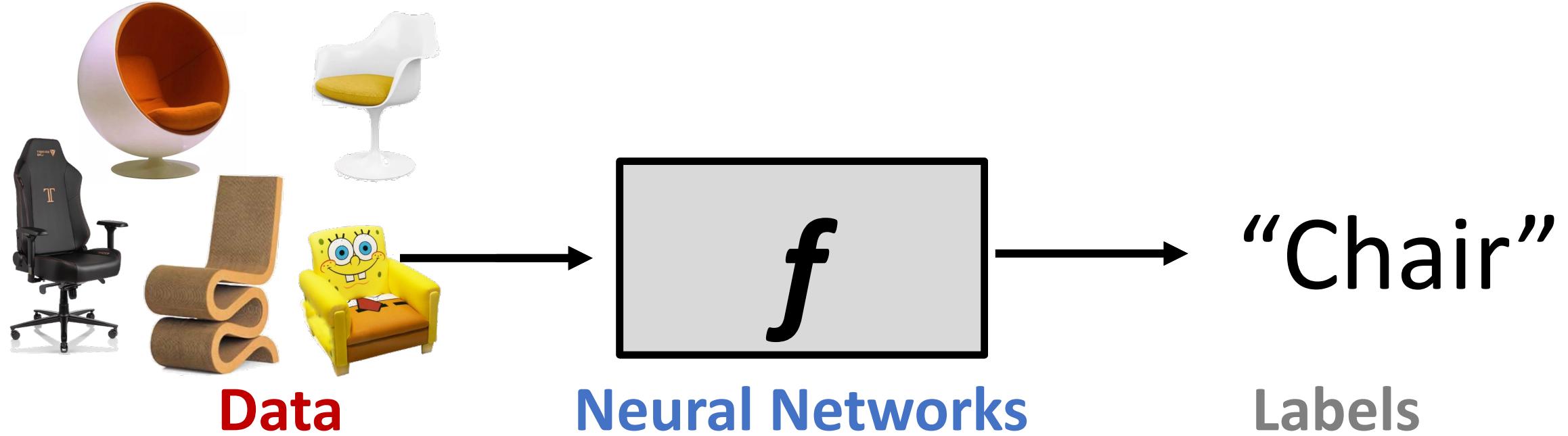
Focused soley on supervised learning so far...



Supervised Learning



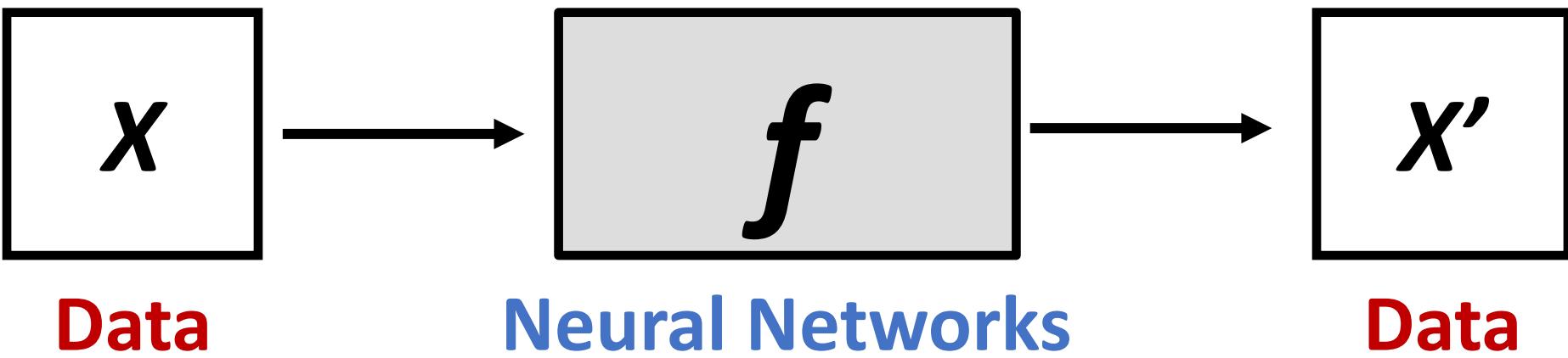
Limitations of Supervised Learning



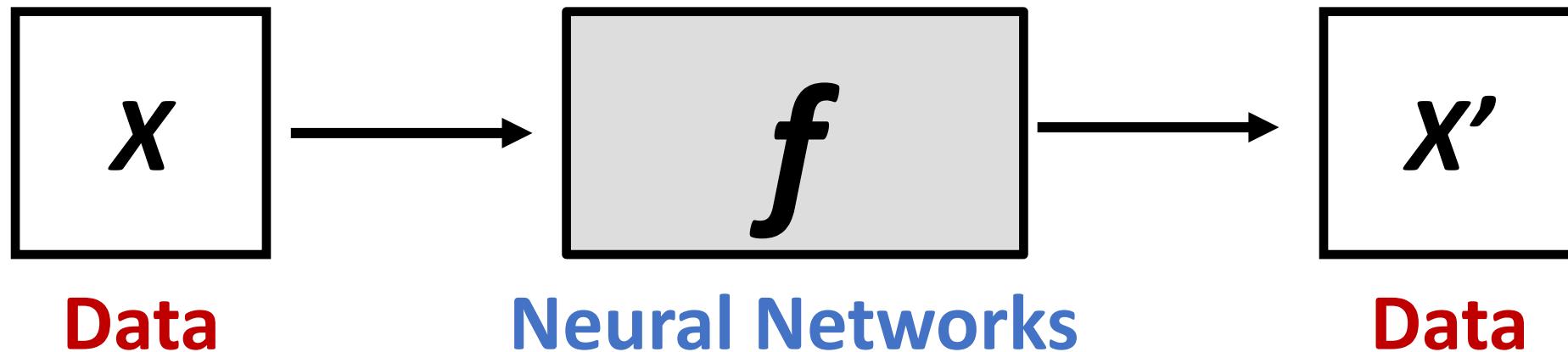
Network might have to cheat:

relying on spurious correlation / memorizing => poor generalization / not robust

Self-Supervised Learning



Self-Supervised Learning



Instead of asking the machine what an object is, ask for:

- Similarity
- Correspondence
- Association
- Prediction
- Reconstruction

Self-Supervised Learning

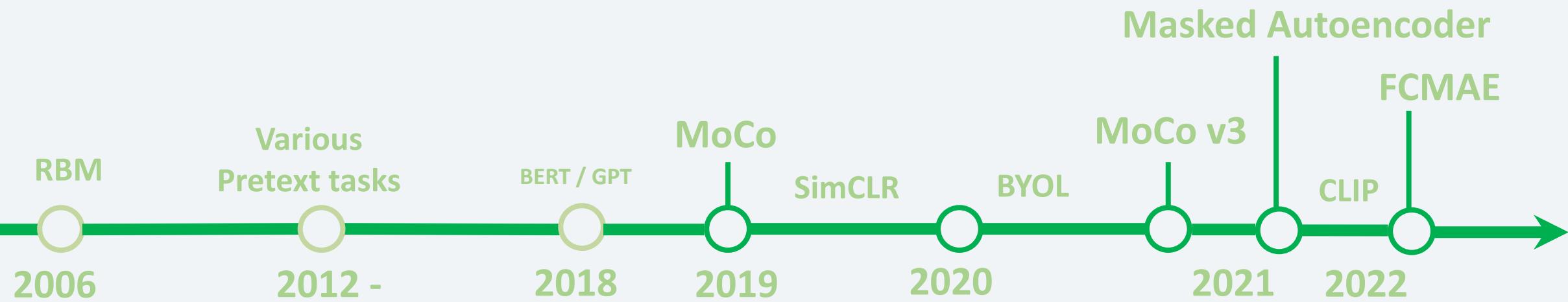


Goal:

To build background knowledge
and approximate a form of common sense in AI
systems.

Architecture / Objective / Data

2. Training objectives beyond supervised classification: Are labels necessary?



Self-supervised Learning



Jigsaw Puzzle

[Noroozi and Favaro, 2016]

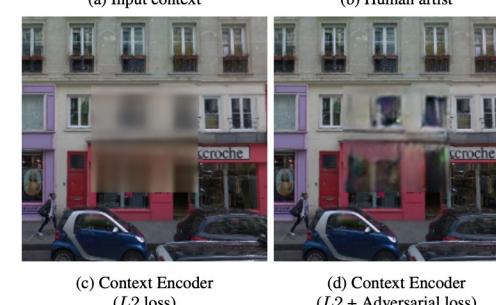
Various
Pretext tasks

2006

2012 -



Colorization
[Zhang et al., 2016]



ContextEncoder
[Pathak et al., 2016]

ContextEncoder

[Pathak et al., 2016]

Pretraining Method	Supervision	Pretraining time	Classification	Detection	Segmentation
ImageNet [26]	1000 class labels	3 days	78.2%	56.8%	48.0%
Random Gaussian	initialization	< 1 minute	53.3%	43.4%	19.8%
Autoencoder	-	14 hours	53.8%	41.9%	25.2%
Agrawal <i>et al.</i> [1]	egomotion	10 hours	52.9%	41.8%	-
Wang <i>et al.</i> [39]	motion	1 week	58.7%	47.4%	-
Doersch <i>et al.</i> [7]	relative context	4 weeks	55.3%	46.6%	-
Ours	context	14 hours	56.5%	44.5%	30.0%

-10% to - 20% compared to ImageNet supervised pre-training
Labels seemed necessary...

Various
Pretext tasks



Contrastive Learning

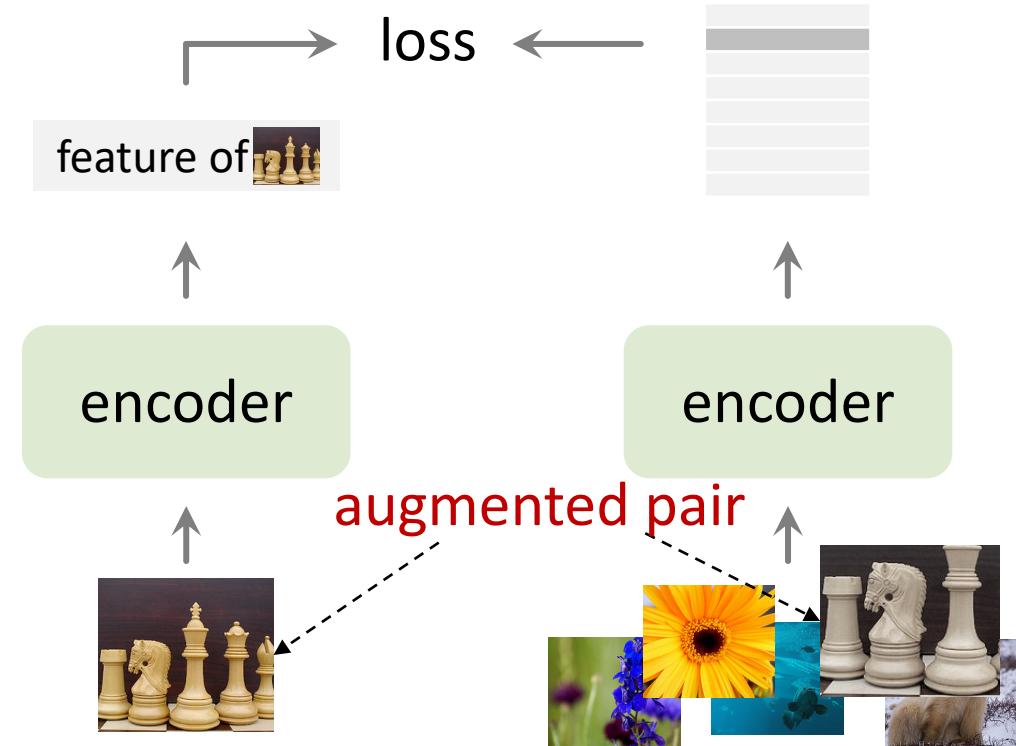
[Chopra *et al.* CVPR 2005]

[Hadsell *et al.* CVPR 2006]

[Wu *et al.* CVPR 2018]

...

Goal: similarity learning



Contrastive Learning

[Chopra *et al.* CVPR 2005]

[Hadsell *et al.* CVPR 2006]

[Wu *et al.* CVPR 2018]

...



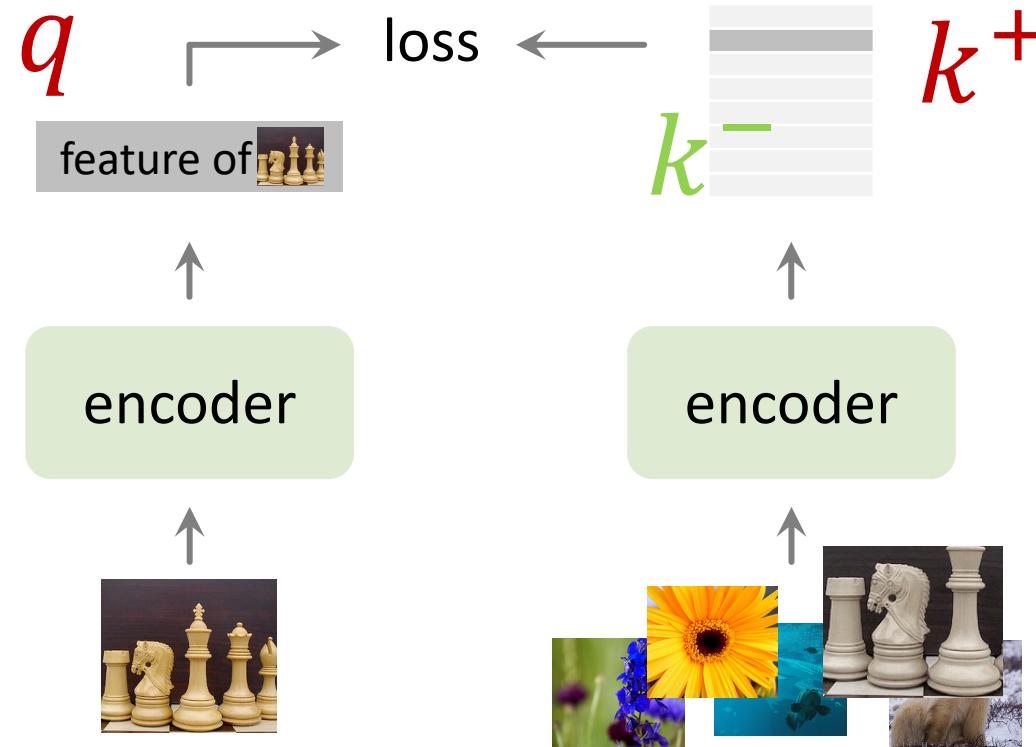
Contrastive Learning

[Chopra *et al.* CVPR 2005]

[Hadsell *et al.* CVPR 2006]

[Wu *et al.* CVPR 2018]

...



Contrastive Learning

[Chopra *et al.* CVPR 2005]

[Hadsell *et al.* CVPR 2006]

[Wu *et al.* CVPR 2018]

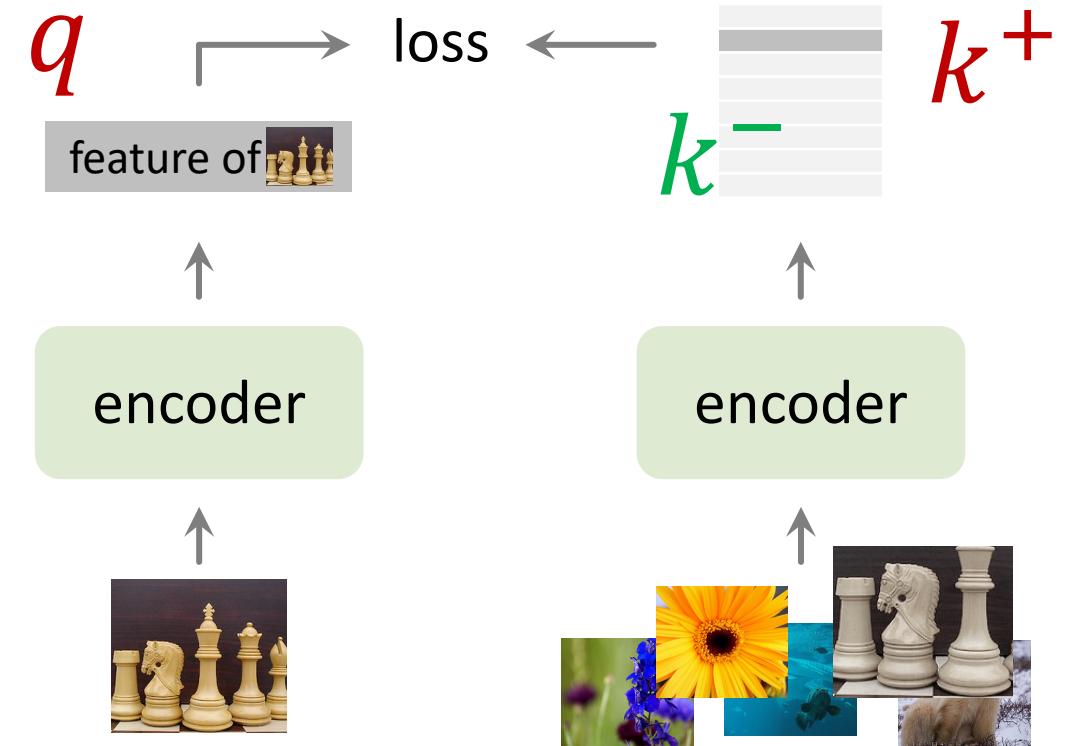
...

InfoNCE Loss:

$$\mathcal{L}_q = -\log \frac{\exp(q \cdot k^+ / \tau)}{\exp(q \cdot k^+ / \tau) + \sum_{k^-} \exp(q \cdot k^- / \tau)}.$$

To learn features that

- attract **similar** samples (q and k^+)
- dispel **dissimilar** samples (q and many k^-).



Momentum Contrast

[He, Fan, Wu, Xie, Girshick. CVPR 2020]

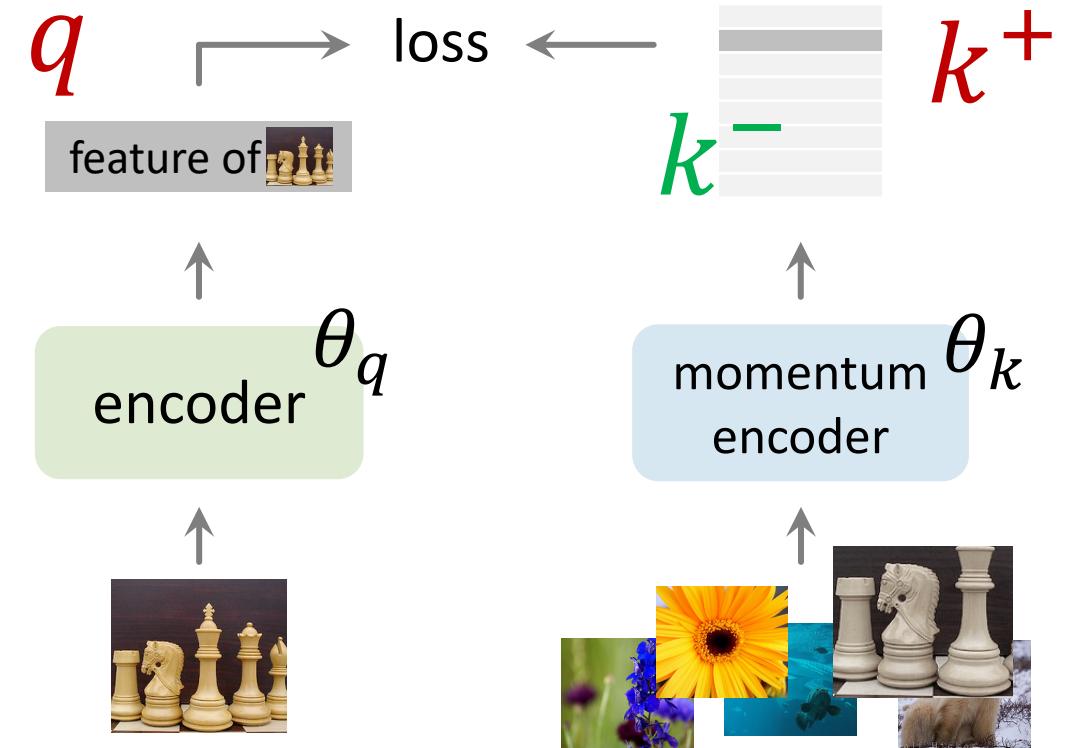
Momentum Encoder:

$$\theta_k := m \cdot \theta_k + (1 - m) \cdot \theta_q$$

slowly update, e.g.
 $m = 0.999$

Enables:

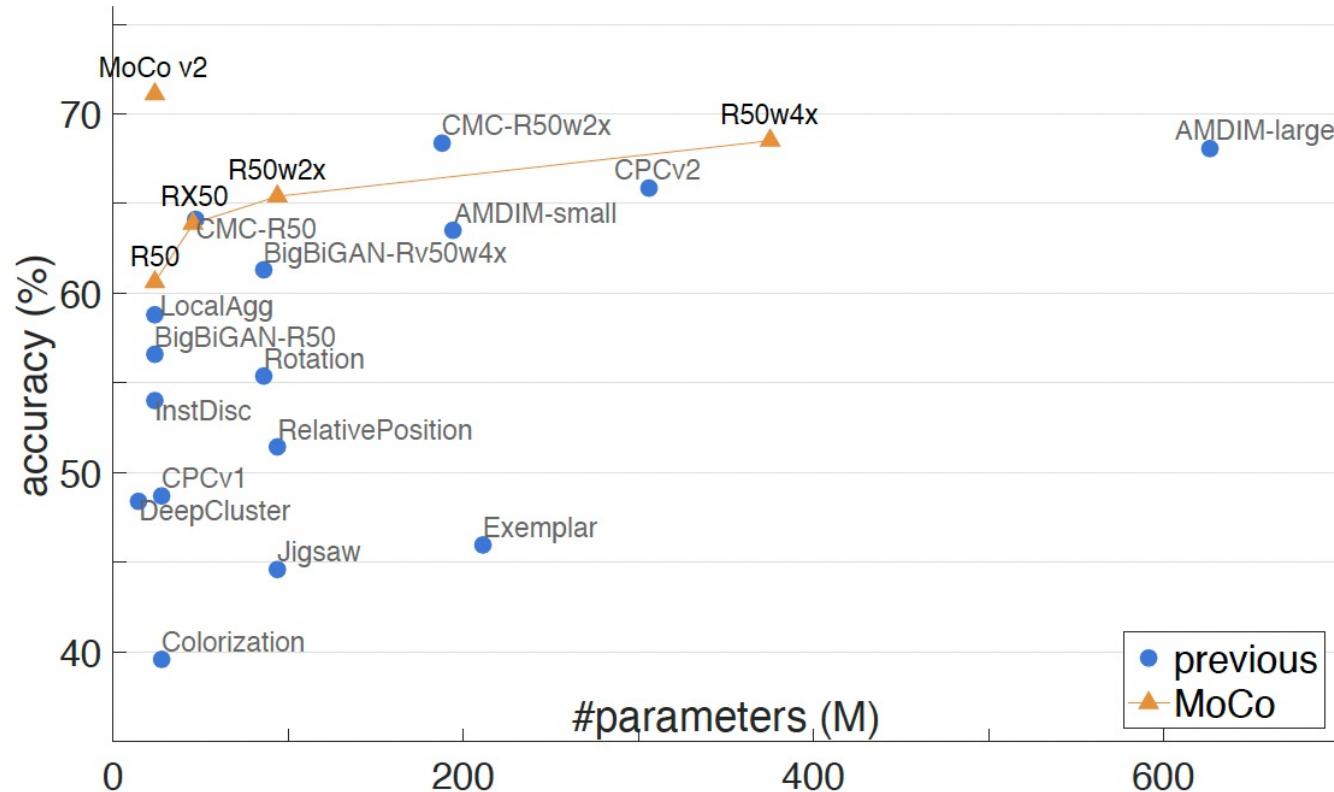
- Large feature queues
- Consistent features



Momentum Contrast

[He, Fan, Wu, **Xie**, Girshick. CVPR 2020]

ImageNet
Linear evaluation



pre-train	AP ₅₀	AP	AP ₇₅
random init.	64.4	37.9	38.6
super. IN-1M	81.4	54.0	59.1
MoCo IN-1M	81.1 (-0.3)	54.6 (+0.6)	59.9 (+0.8)
MoCo IG-1B	81.6 (+0.2)	55.5 (+1.5)	61.2 (+2.1)

(a) Faster R-CNN, R50-dilated-C5

pre-train	AP ₅₀	AP	AP ₇₅
random init.	60.2	33.8	33.1
super. IN-1M	81.3	53.5	58.8
MoCo IN-1M	81.5 (+0.2)	55.9 (+2.4)	62.6 (+3.8)
MoCo IG-1B	82.2 (+0.9)	57.2 (+3.7)	63.7 (+4.9)

(b) Faster R-CNN, R50-C4

AP ^{bb}	AP ₅₀ ^{bb}	AP ₇₅ ^{bb}	AP ^{mk}	AP ₅₀ ^{mk}	AP ₇₅ ^{mk}
36.7	56.7	40.0	33.7	53.8	35.9
40.6	61.3	44.4	36.8	58.1	39.5
40.8 (+0.2)	61.6 (+0.3)	44.7 (+0.3)	36.9 (+0.1)	58.4 (+0.3)	39.7 (+0.2)
41.1 (+0.5)	61.8 (+0.5)	45.1 (+0.7)	37.4 (+0.6)	59.1 (+1.0)	40.2 (+0.7)

(b) Mask R-CNN, R50-FPN, 2× schedule

AP ^{bb}	AP ₅₀ ^{bb}	AP ₇₅ ^{bb}	AP ^{mk}	AP ₅₀ ^{mk}	AP ₇₅ ^{mk}
35.6	54.6	38.2	31.4	51.5	33.5
40.0	59.9	43.1	34.7	56.5	36.9
40.7 (+0.7)	60.5 (+0.6)	44.1 (+1.0)	35.4 (+0.7)	57.3 (+0.8)	37.6 (+0.7)
41.1 (+1.1)	60.7 (+0.8)	44.8 (+1.7)	35.6 (+0.9)	57.4 (+0.9)	38.1 (+1.2)

(d) Mask R-CNN, R50-C4, 2× schedule

VOC 07+12 Detection
surpass, +4.9 AP₇₅

COCO Detection
COCO Instance seg.
surpass

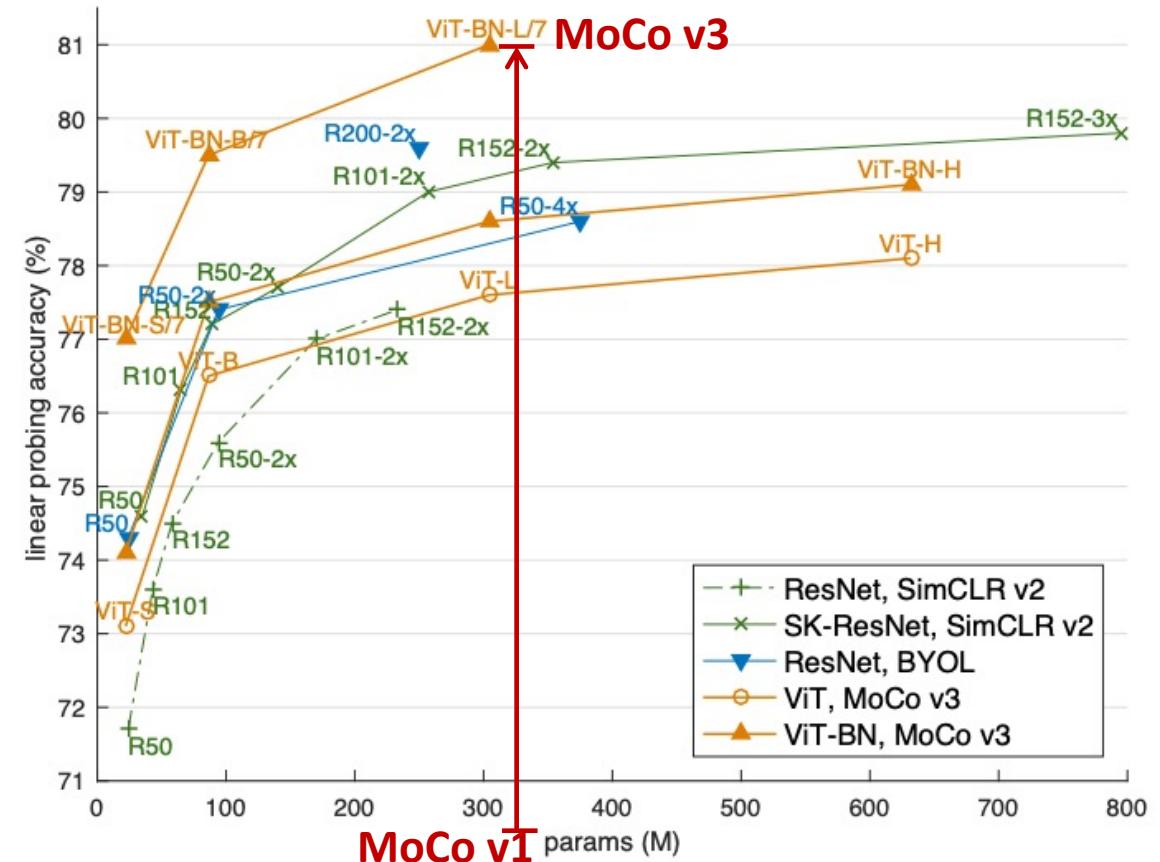
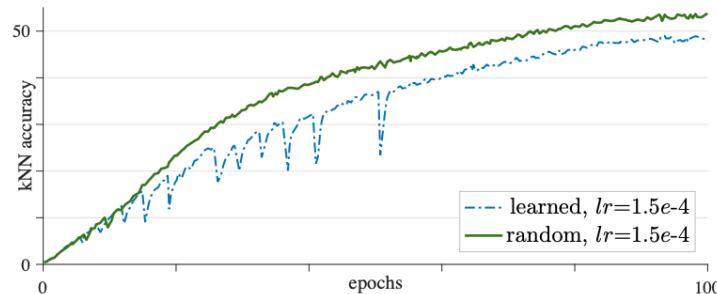
Self-supervised pre-training: surpass supervised counterparts in 7 vision tasks



MoCo v1 to v3

[Chen*, Xie*, He, ICCV 2021]

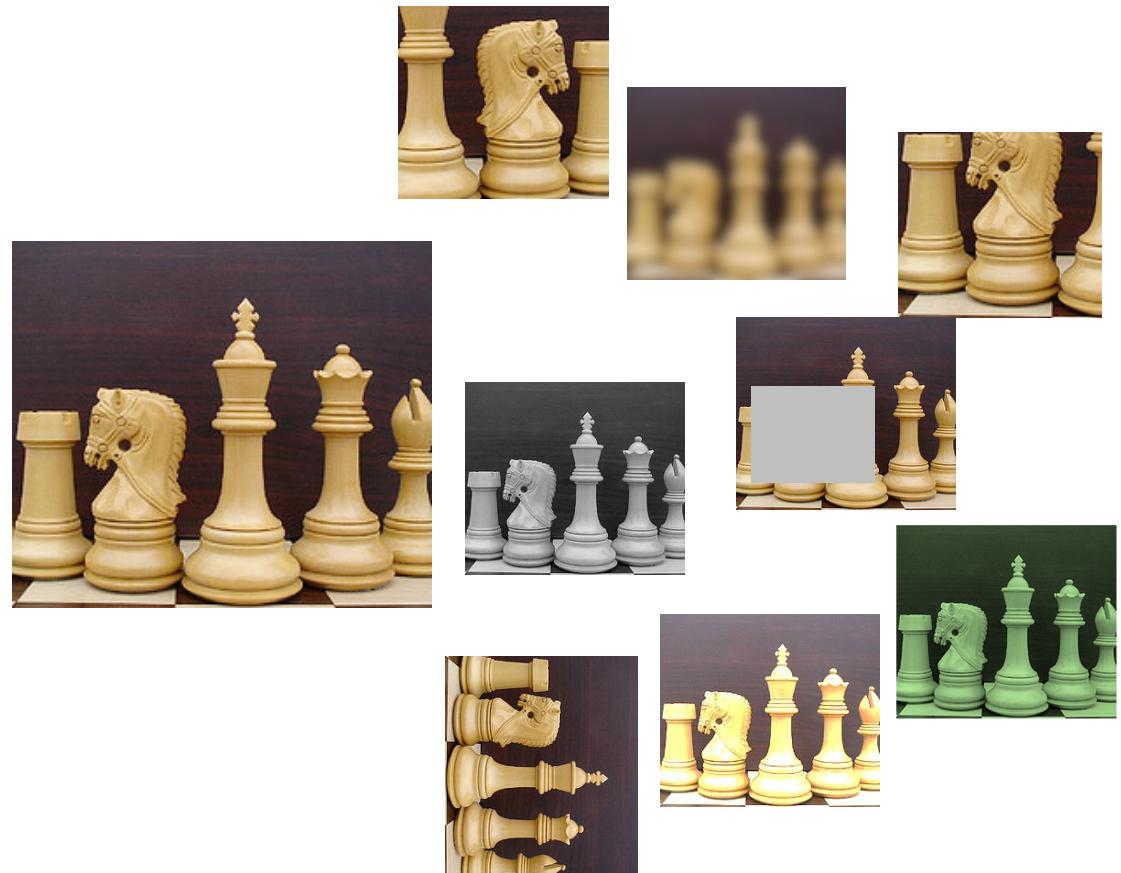
- Using bigger models
- Better training techniques
- Stabilize training with *random patch projection*



Contrastive Learning

Rely on:

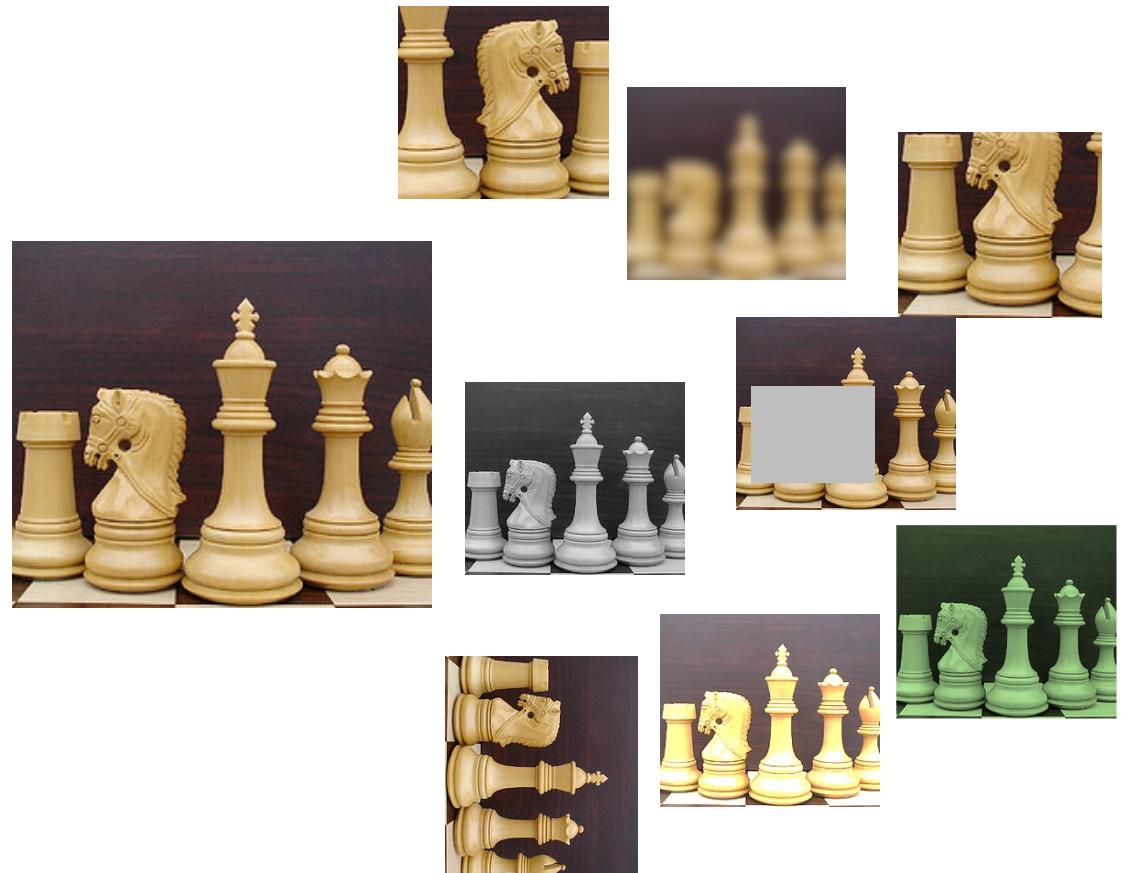
- hand-designed *data augmentations*
- “Too easy to solve” otherwise



Contrastive Learning

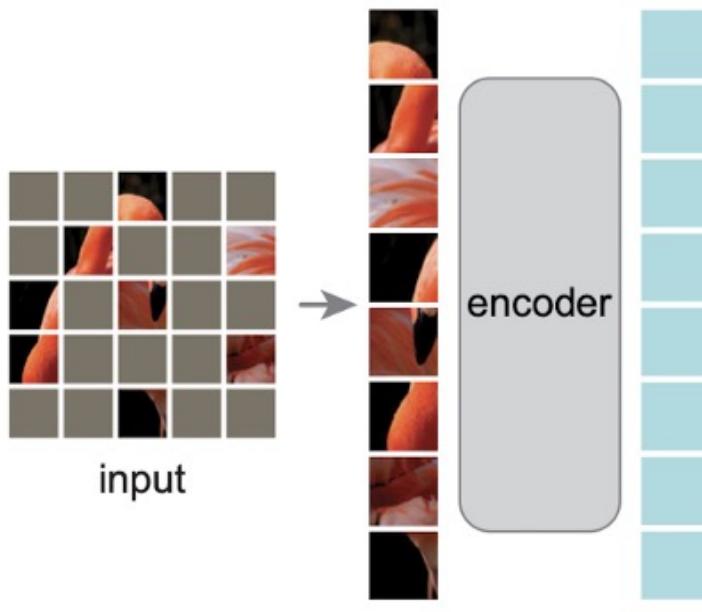
Rely on:

- negatives to prevent collapse
- hand-designed *data augmentations*
“Too easy to solve” otherwise
- What are other possible tasks?
 - “Classical methods” revisited
 - Clustering? 😕
 - Autoencoder? 😕
 - Denoising autoencoder? 🤔



Masked Autoencoders (MAE)

[He, Chen, Xie, Li, Dollár, Girshick, CVPR 2022]

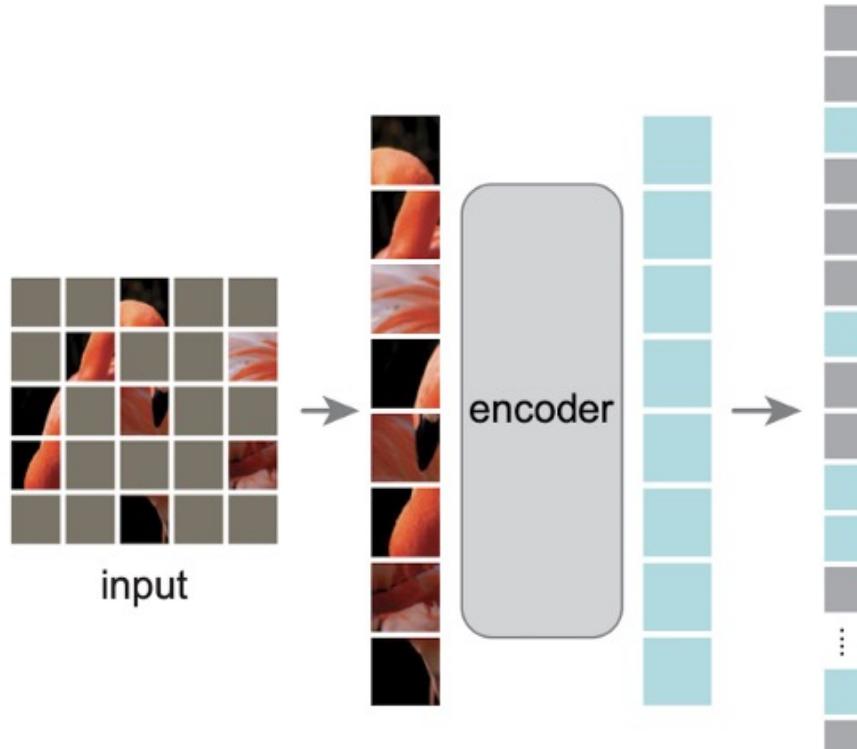


encode visible patches



Masked Autoencoders (MAE)

[He, Chen, Xie, Li, Dollár, Girshick, CVPR 2022]



add “mask tokens”



MAE

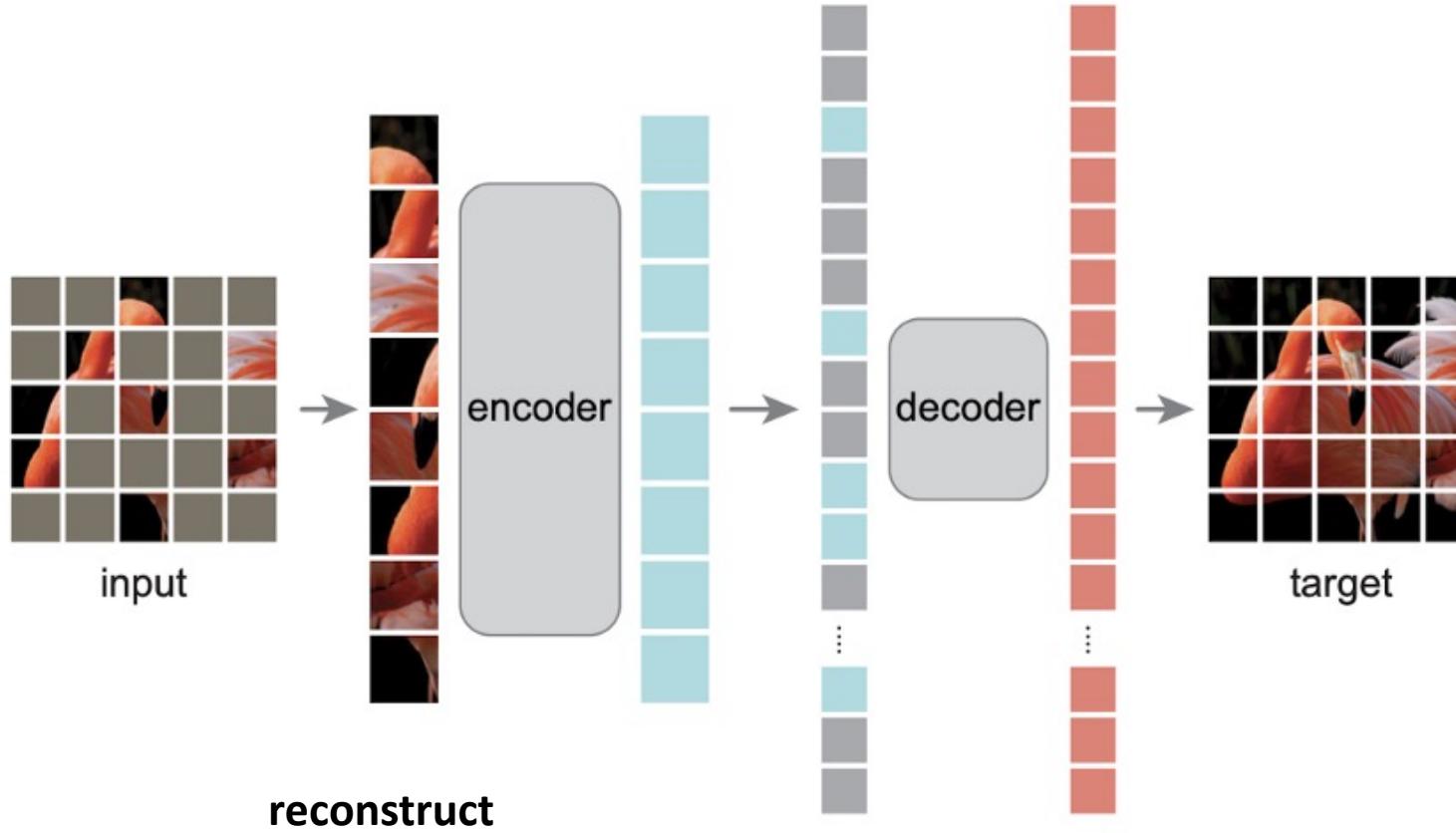
2021

2018

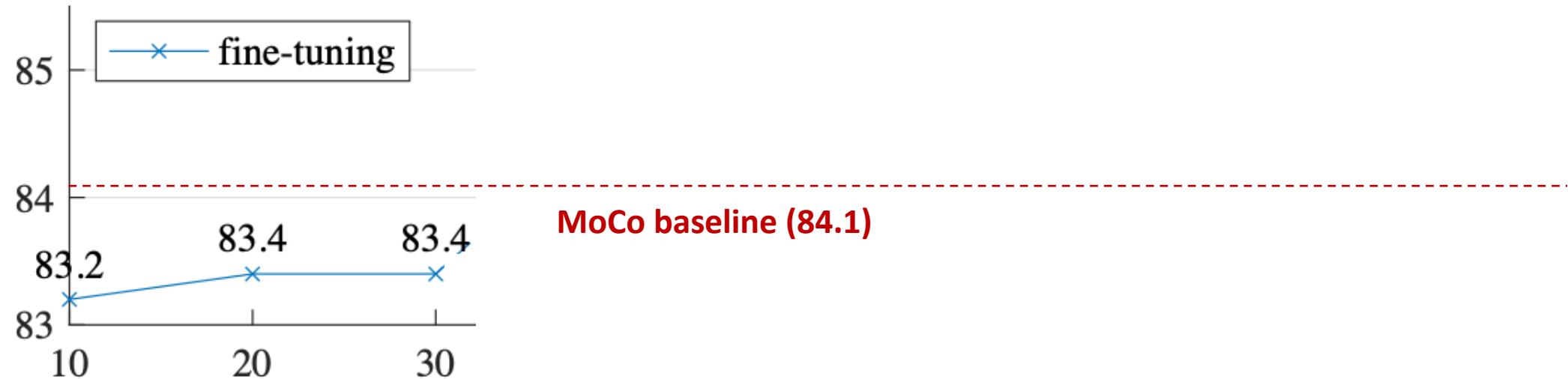
2019

Masked Autoencoders (MAE)

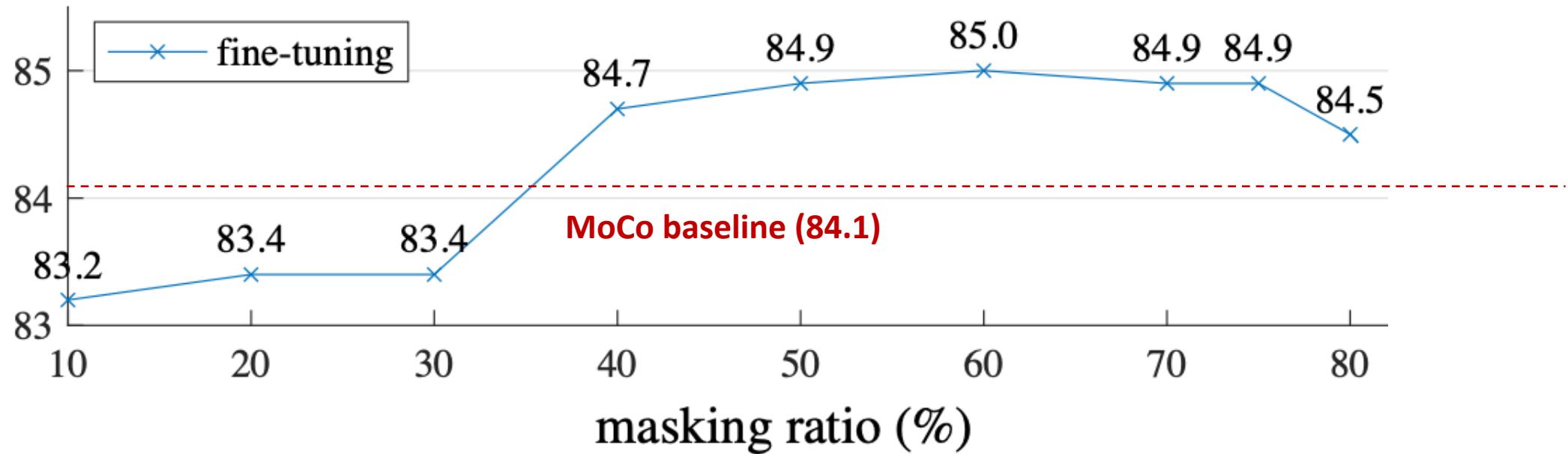
[He, Chen, Xie, Li, Dollár, Girshick, CVPR 2022]



Masking Ratio



Masking Ratio



Visualization



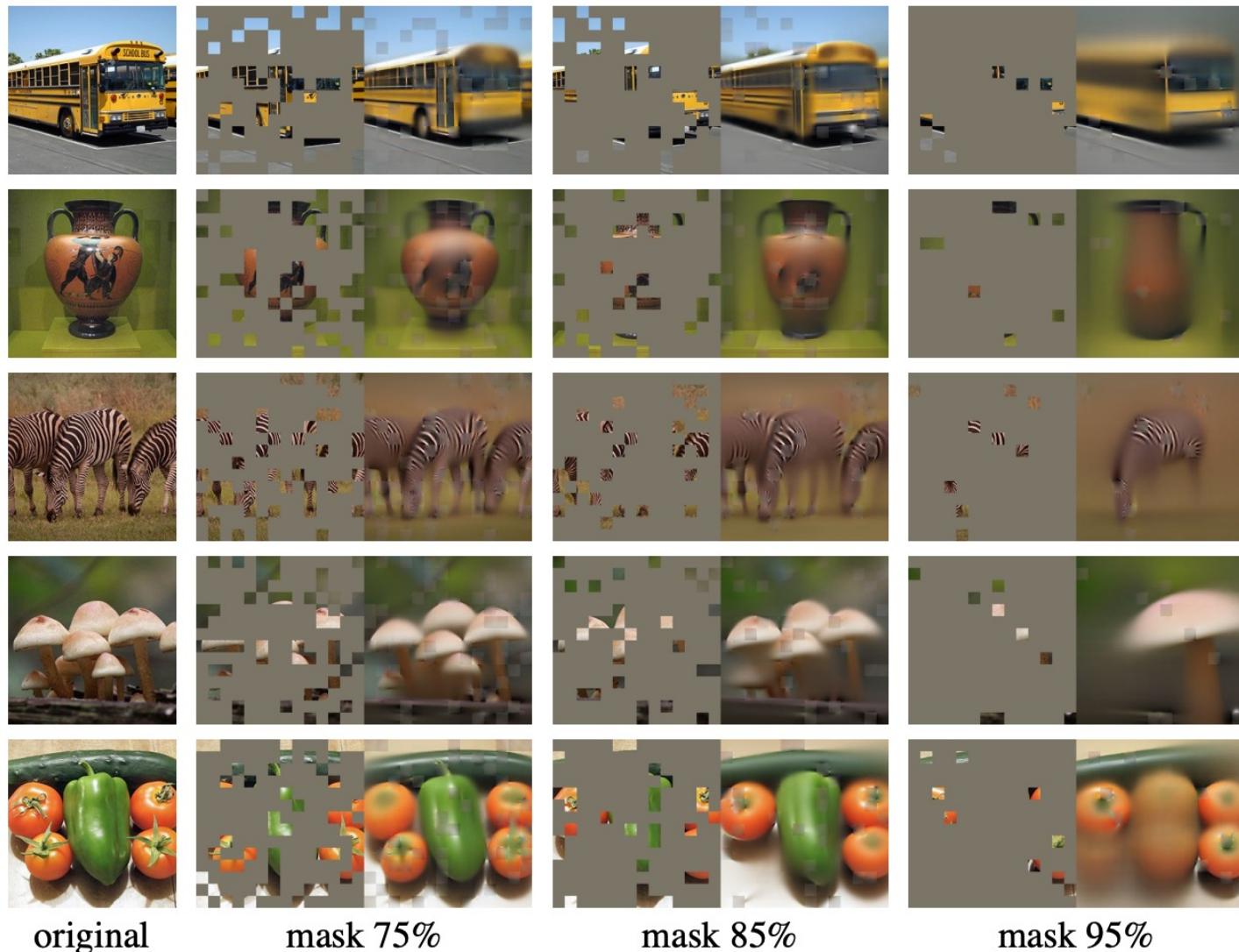
Why high masking ratio

- Eliminates redundancy
- Cannot be easily solved by extrapolating from visible neighboring patches
- Practical benefits:
significant speedup!

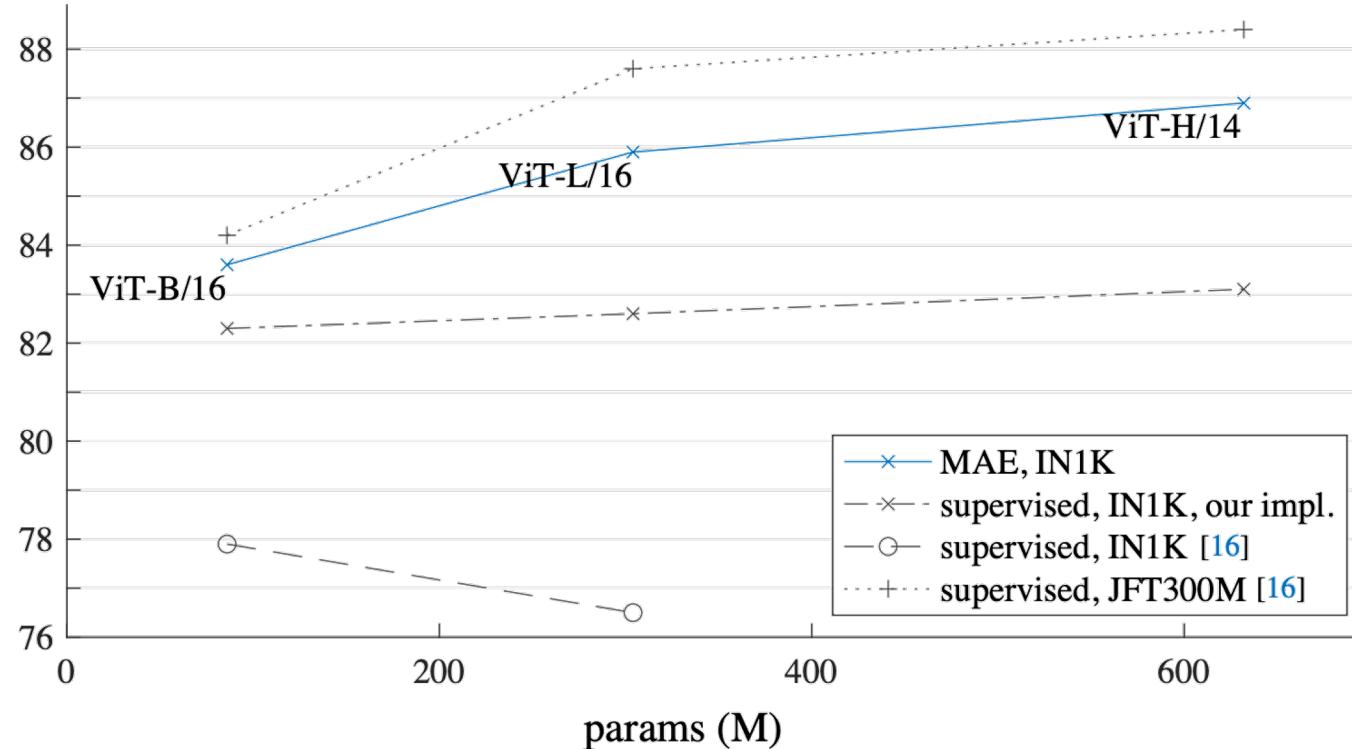


Generalization

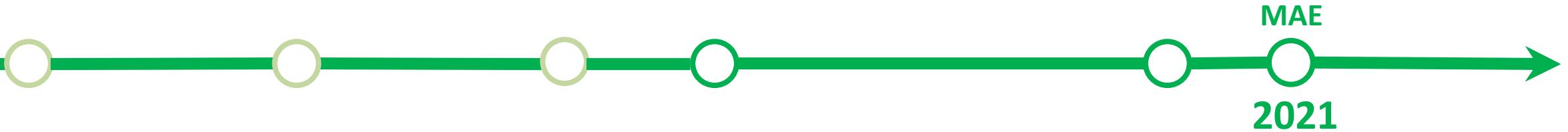
- Samples from the *validation* set
 - Pre-trained with 75% masking ratio
 - Applied on higher masking ratios
- Predictions differ plausibly
 - Representations are nontrivial
 - Model can generalize



Results



- New SOTA in ImageNet-1K (no external data): 87.8%



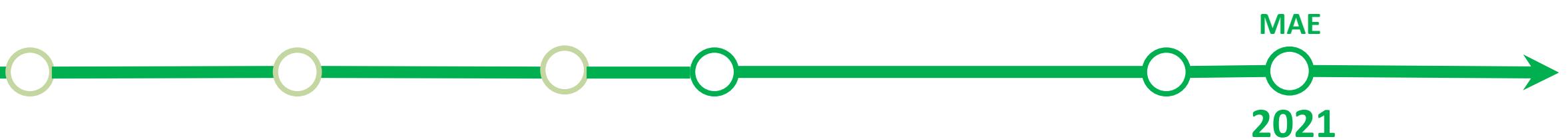
Results

method	pre-train data	AP ^{box}		AP ^{mask}		method	pre-train data	ViT-B	ViT-L
		ViT-B	ViT-L	ViT-B	ViT-L				
supervised	IN1K w/ labels	47.9	49.3	42.9	43.9	supervised	IN1K w/ labels	47.4	49.9
MoCo v3	IN1K	47.9	49.3	42.7	44.0	MoCo v3	IN1K	47.3	49.1
BEiT	IN1K+DALLE	49.8	53.3	44.4	47.1	BEiT	IN1K+DALLE	47.1	53.3
MAE	IN1K	50.3	53.3	44.9	47.2	MAE	IN1K	48.1	53.6

- Transfer Learning on detection / segmentation

dataset	ViT-B	ViT-L	ViT-H	ViT-H ₄₄₈	prev best
IN-Corruption ↓ [27]	51.7	41.8	33.8	36.8	42.5 [32]
IN-Adversarial [28]	35.9	57.1	68.2	76.7	35.8 [41]
IN-Rendition [26]	48.3	59.9	64.4	66.5	48.7 [41]
IN-Sketch [60]	34.5	45.3	49.6	50.9	36.0 [41]

- Robust evaluations: +40% previous SOTA



Rethinking model robustness

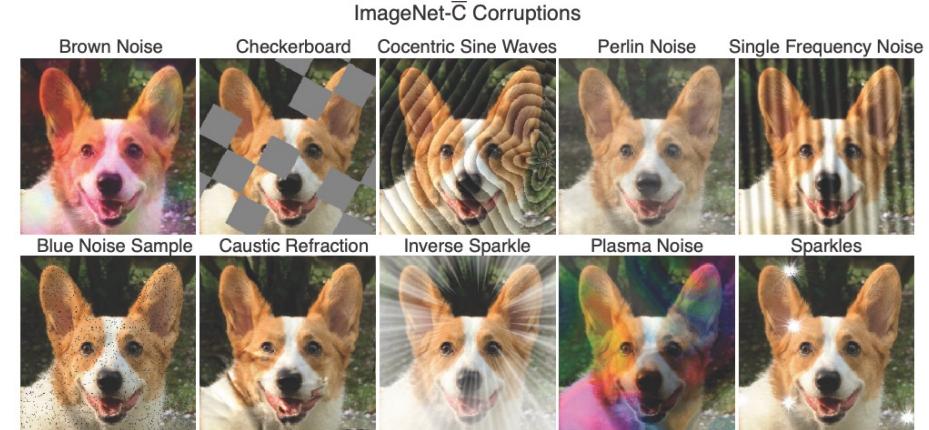
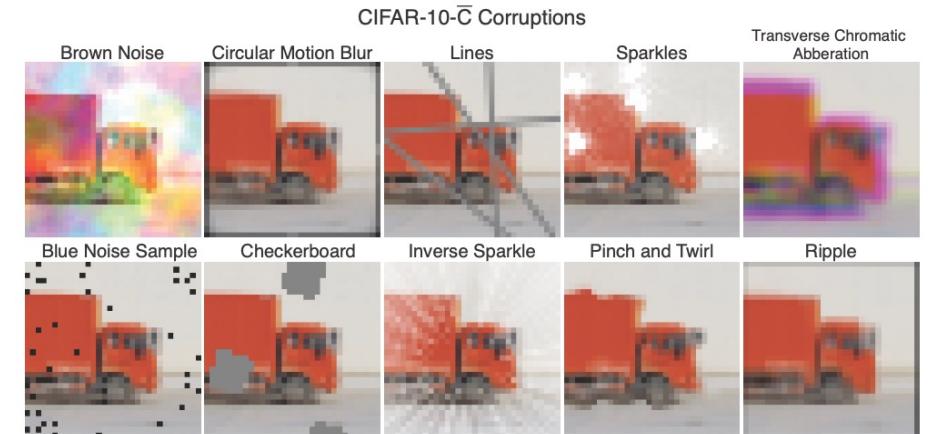
Masked Autoencoder

dataset	ViT-B	ViT-L	ViT-H	ViT-H ₄₄₈	prev best
IN-Corruption ↓ [27]	51.7	41.8	33.8	36.8	42.5 [32]
IN-Adversarial [28]	35.9	57.1	68.2	76.7	35.8 [41]
IN-Rendition [26]	48.3	59.9	64.4	66.5	48.7 [41]
IN-Sketch [60]	34.5	45.3	49.6	50.9	36.0 [41]

Scaling brings +40% in accuracy

ConvNeXt

Model	Data/Size	FLOPs / Params	Clean	C (↓)	Ā (↓)	A	R	SK
ResNet-50	1K/224 ²	4.1 / 25.6	76.1	76.7	57.7	0.0	36.1	24.1
Swin-T [42]	1K/224 ²	4.5 / 28.3	81.2	62.0	-	21.6	41.3	29.1
RVT-S* [44]	1K/224 ²	4.7 / 23.3	81.9	49.4	37.5	25.7	47.7	34.7
ConvNeXt-T	1K/224 ²	4.5 / 28.6	82.1	53.2	40.0	24.2	47.2	33.8
Swin-B [42]	1K/224 ²	15.4 / 87.8	83.4	54.4	-	35.8	46.6	32.4
RVT-B* [44]	1K/224 ²	17.7 / 91.8	82.6	46.8	30.8	28.5	48.7	36.0
ConvNeXt-B	1K/224 ²	15.4 / 88.6	83.8	46.8	34.4	36.7	51.3	38.2
ConvNeXt-B	22K/384 ²	45.1 / 88.6	86.8	43.1	30.7	62.3	64.9	51.6
ConvNeXt-L	22K/384 ²	101.0 / 197.8	87.5	40.2	29.9	65.5	66.7	52.8
ConvNeXt-XL	22K/384 ²	179.0 / 350.2	87.8	38.8	27.1	69.3	68.2	55.0

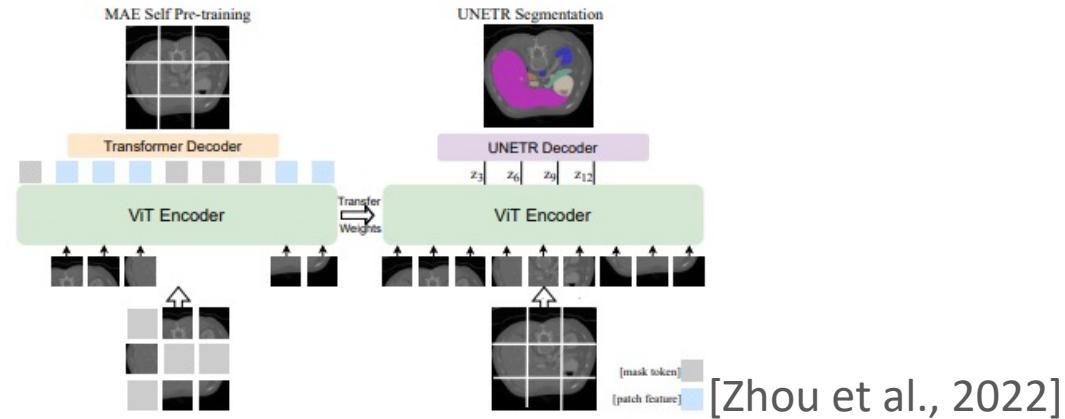


[Mintun, Kirillove, Xie, NeurIPS 2021]

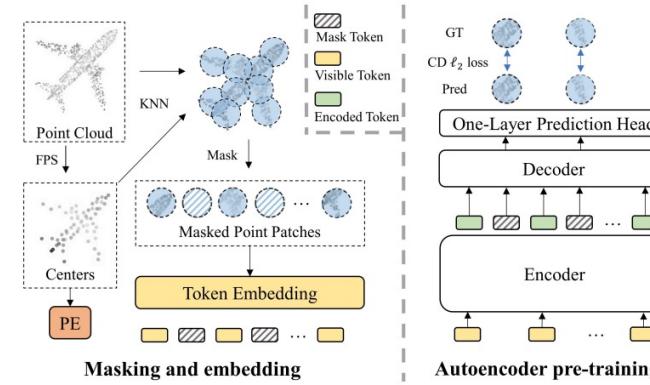
ConvNeXt – similar story

MAE as a general methodology

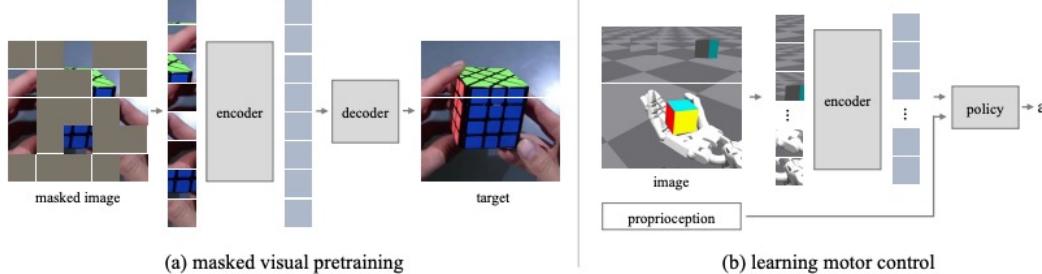
Medical Image Analysis



Point Cloud Classification

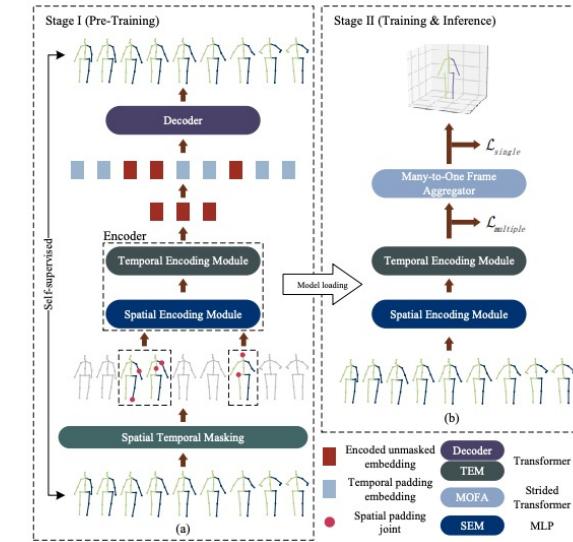


Motor Control

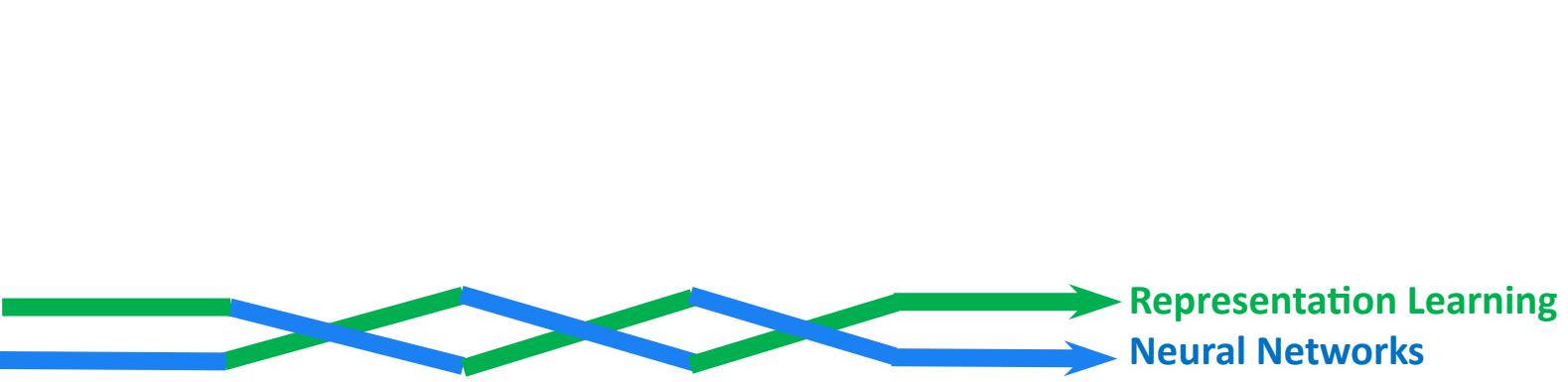


[Xiao, Radosavovic, et al., 2022]

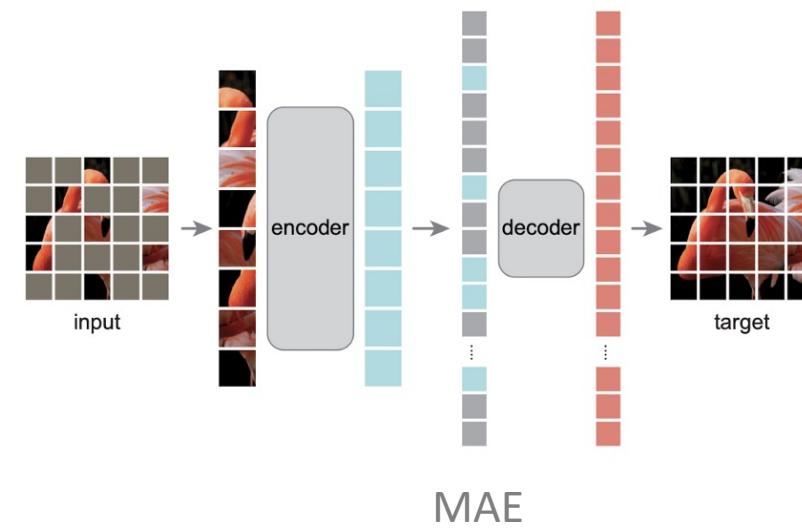
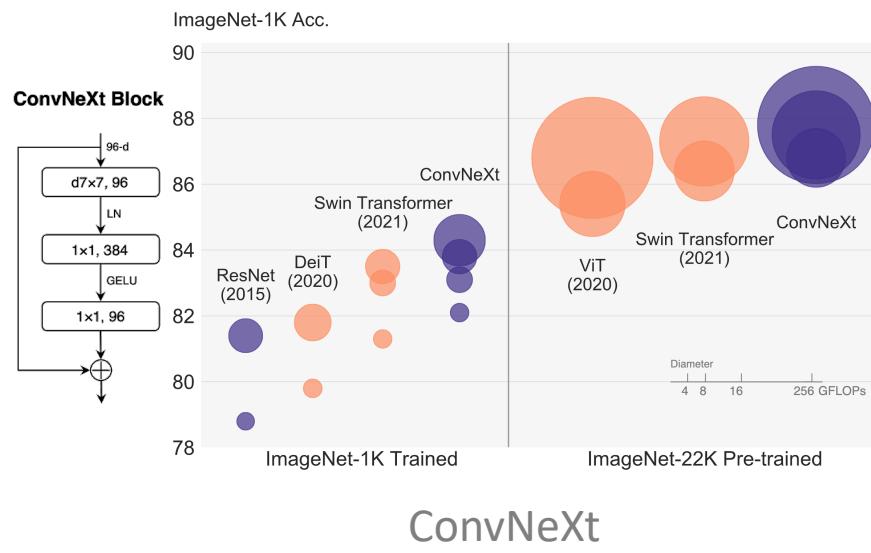
Human Pose Estimation



[Shan et al. 2022]



Common theme – scalable representation learning

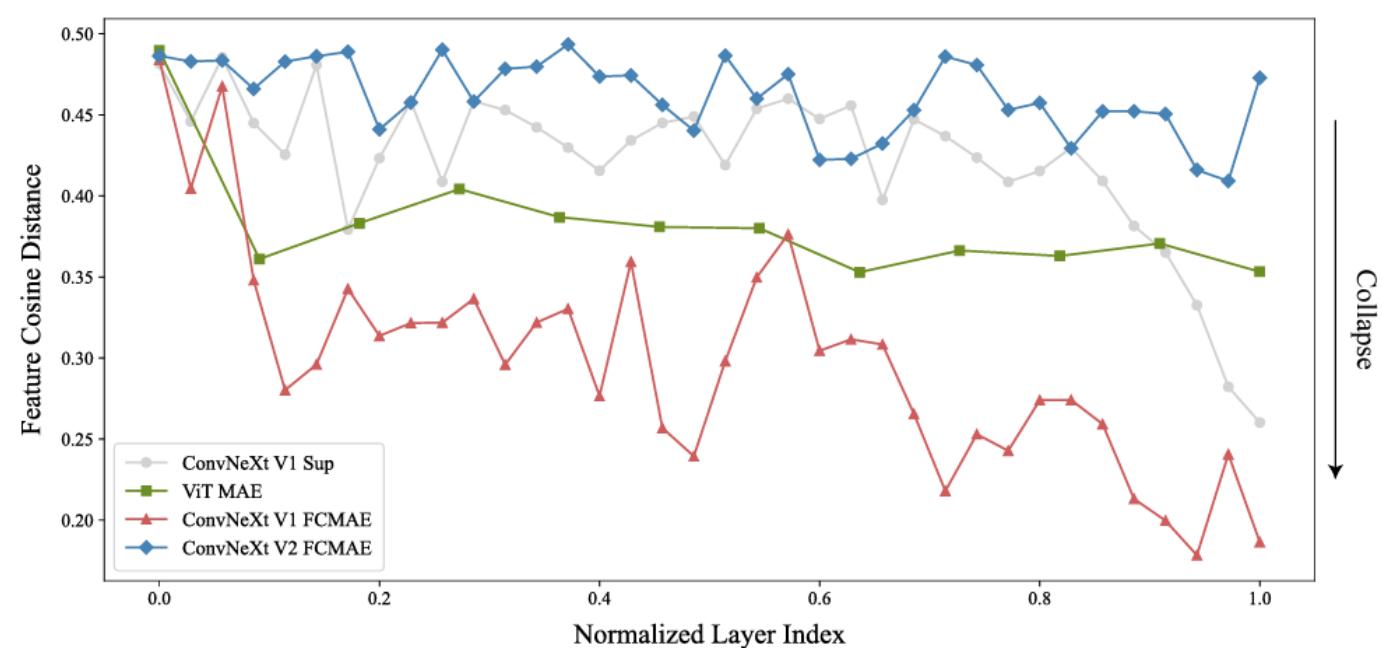


But they are not co-designed...

Can we do MAE with ConvNeXt?

		$\frac{1}{C^2} \sum_i^C \sum_j^C \frac{1 - \cos(X_i, X_j)}{2}$
Sup, 100ep		82.7
Sup, 300ep. [52]		83.8
FCMAE		83.7

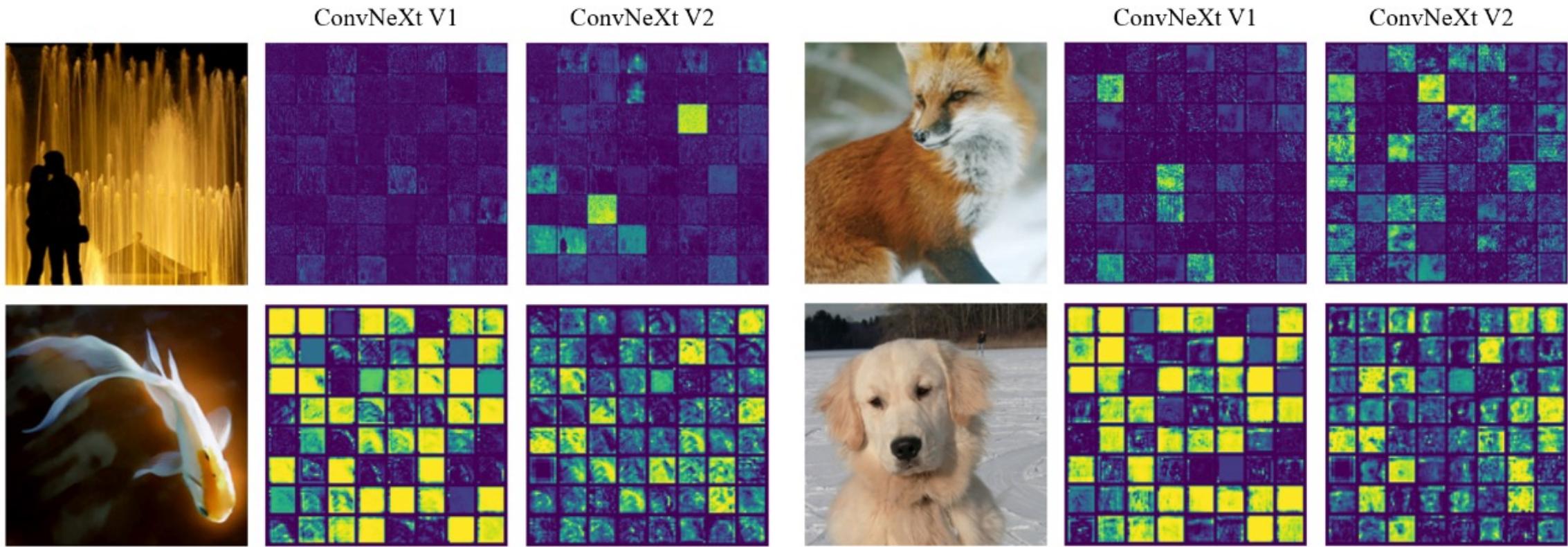
ConvNeXt v1 + MAE doesn't really work (!!)



Potential issue: feature collapse



Feature Collapse



Global Response Normalization (GRN)

Feature Aggregation:

$$\mathcal{G}(X) := X \in \mathcal{R}^{H \times W \times C} \rightarrow gx \in \mathcal{R}^C. \quad (1)$$

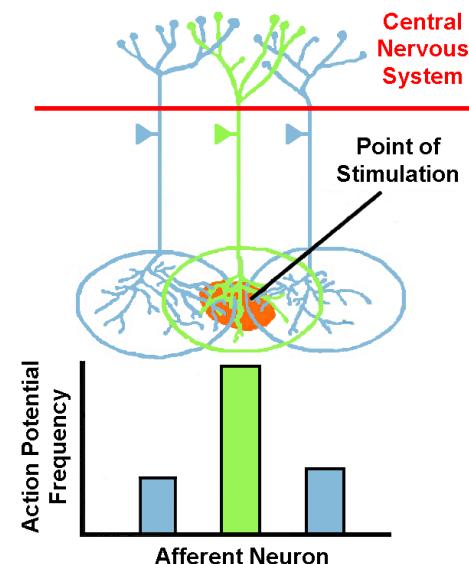
Relative Importance:

$$\mathcal{N}(\|X_i\|) := \|X_i\| \in \mathcal{R} \rightarrow \frac{\|X_i\|}{\sum_{j=1,\dots,C} \|X_j\|} \in \mathcal{R}, \quad (2)$$

Feature Reweighting:

$$X_i = X_i * \mathcal{N}(\mathcal{G}(X)_i) \in \mathcal{R}^{H \times W} \quad (3)$$

Related: Lateral inhibition to sharpen signals to the brain



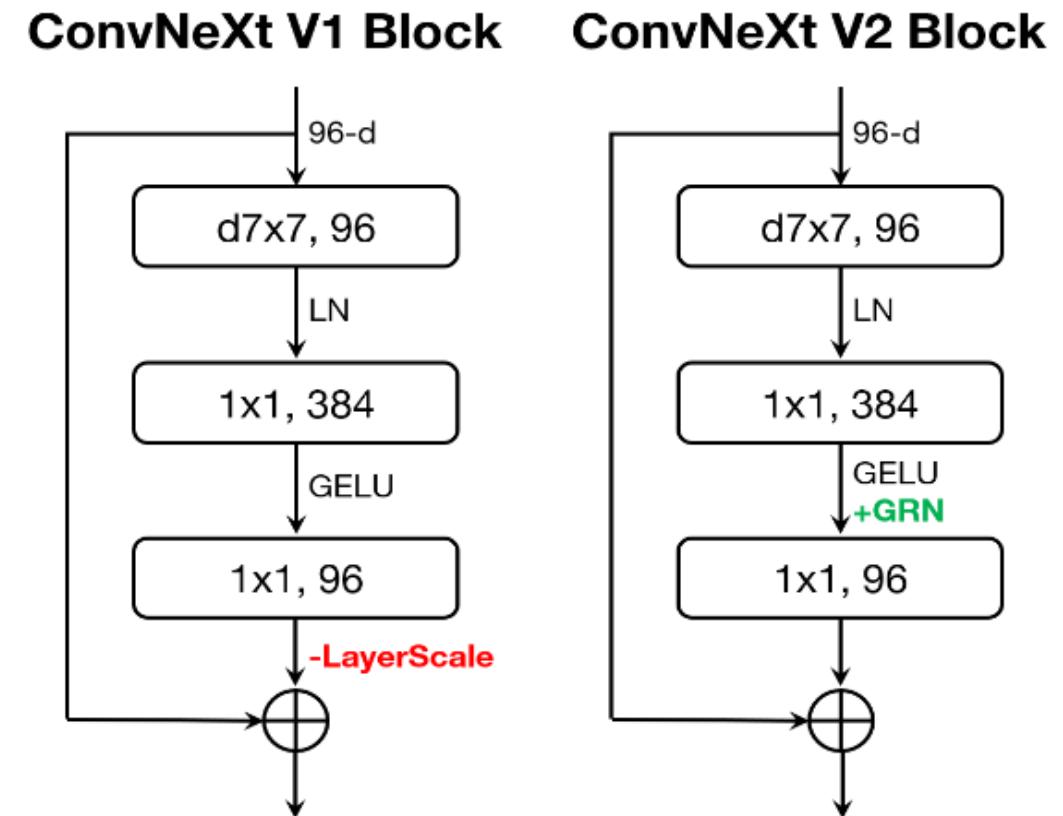
Global Response Normalization (GRN)

Algorithm 1 Pseudocode of GRN in a PyTorch-like style.

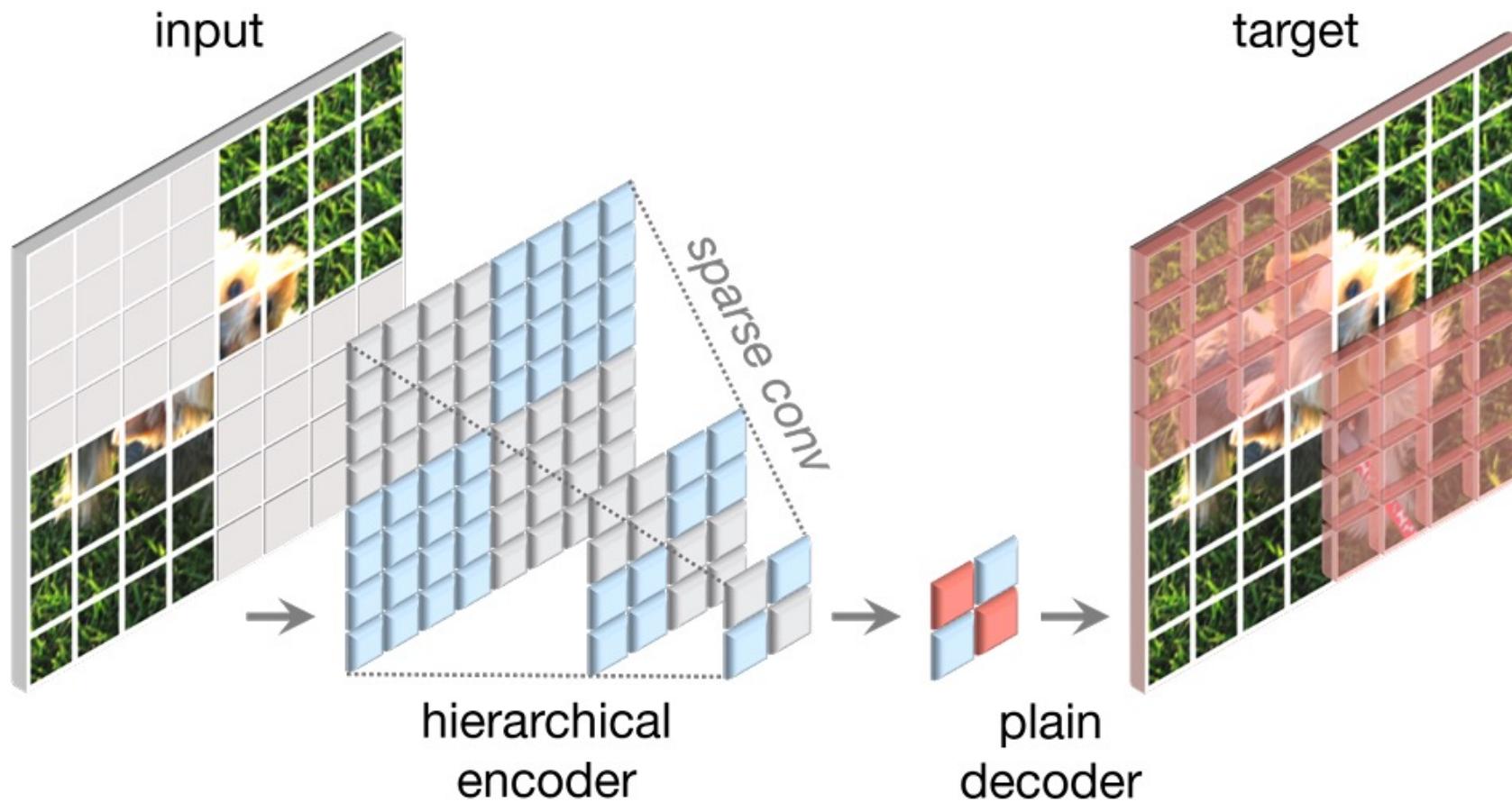
```
# gamma, beta: learnable affine transform parameters
# X: input of shape (N, H, W, C)

gx = torch.norm(X, p=2, dim=(1,2), keepdim=True)
nx = gx / (gx.mean(dim=-1, keepdim=True)+1e-6)
return gamma * (X * nx) + beta + X
```

Only three lines of code
Parameter-free!

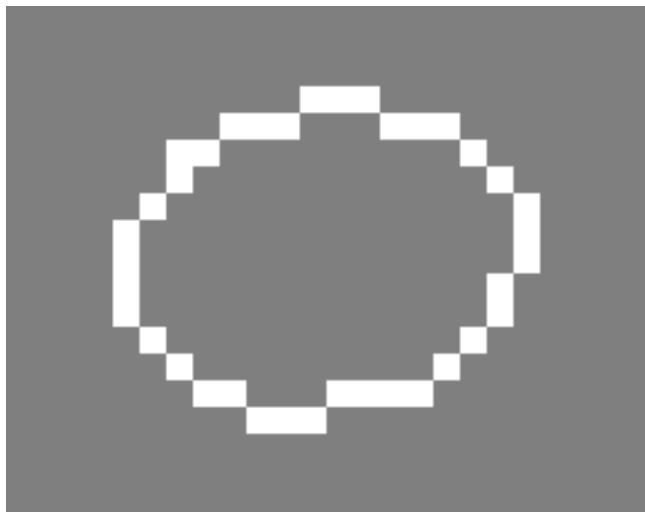


Fully-convolutional MAE

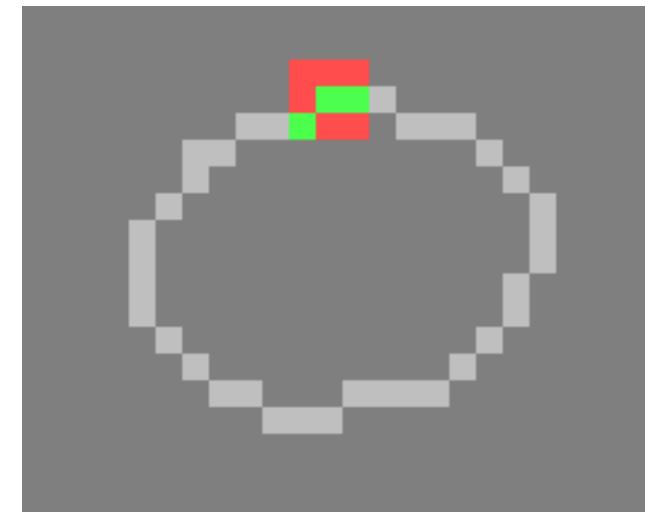


Sparse ConvNet

Non-zero active sites grow rapidly

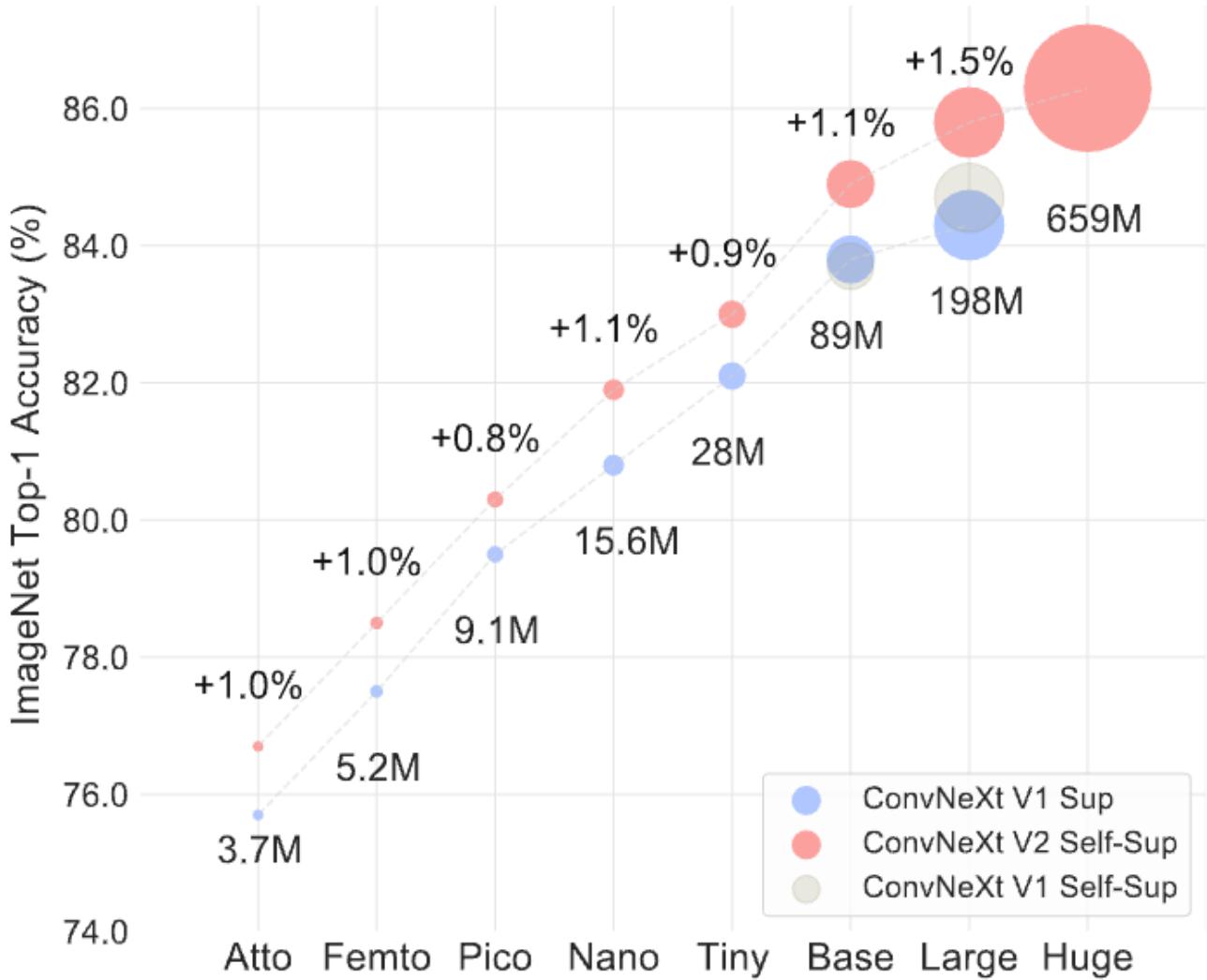


Computation only on a sparse set of “active sites”



Graham, Benjamin, and Laurens van der Maaten. “Submanifold sparse convolutional networks.” *CVPR 2018*

ConvNeXt-v2



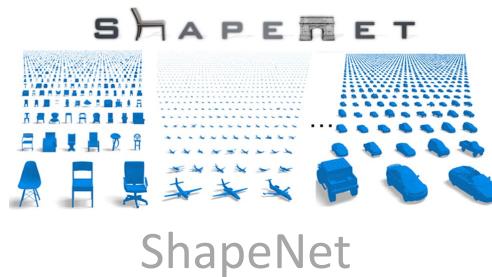
Architecture / Objective / **Data**

3. What data to use for visual pre-training

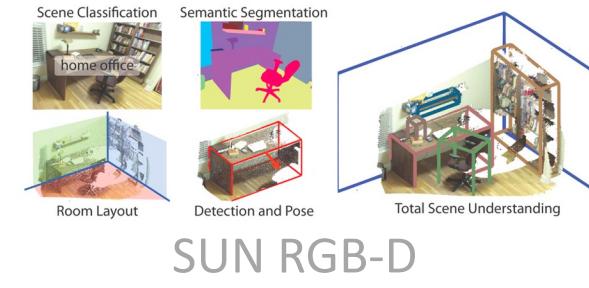
**Multi-Modality
Noisy Labels**



Visual pre-training on 3D point clouds



ShapeNet



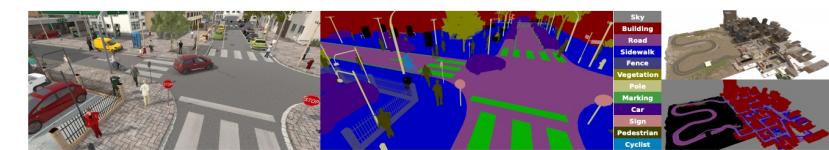
SUN RGB-D



S3DIS



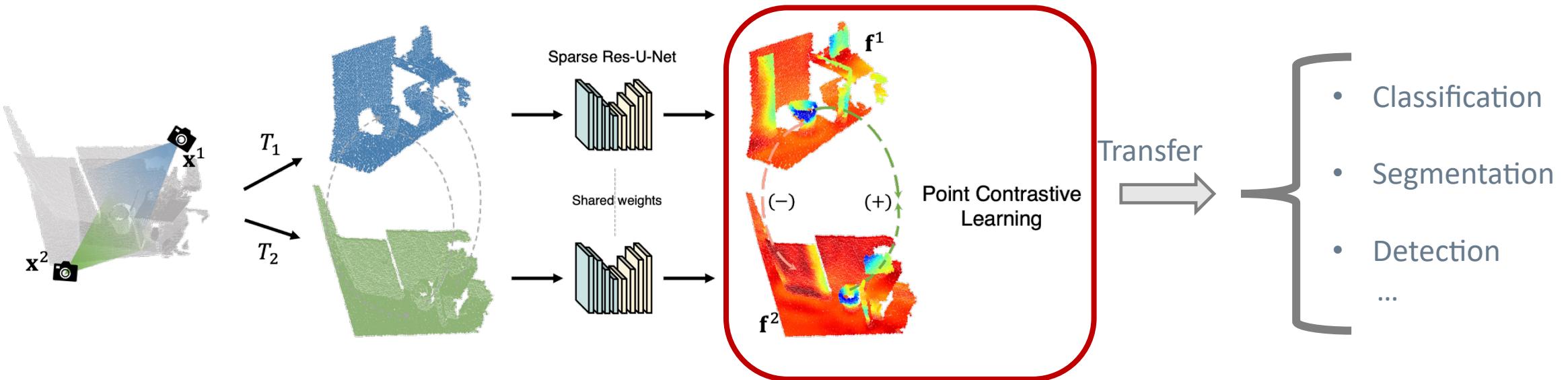
ScanNet



Synthia4D

PointContrast

[Xie, Gu, Guo, Qi, Guibas, Litany, ECCV 2020]
[Hou, Graham, Nießner, Xie, CVPR 2021]



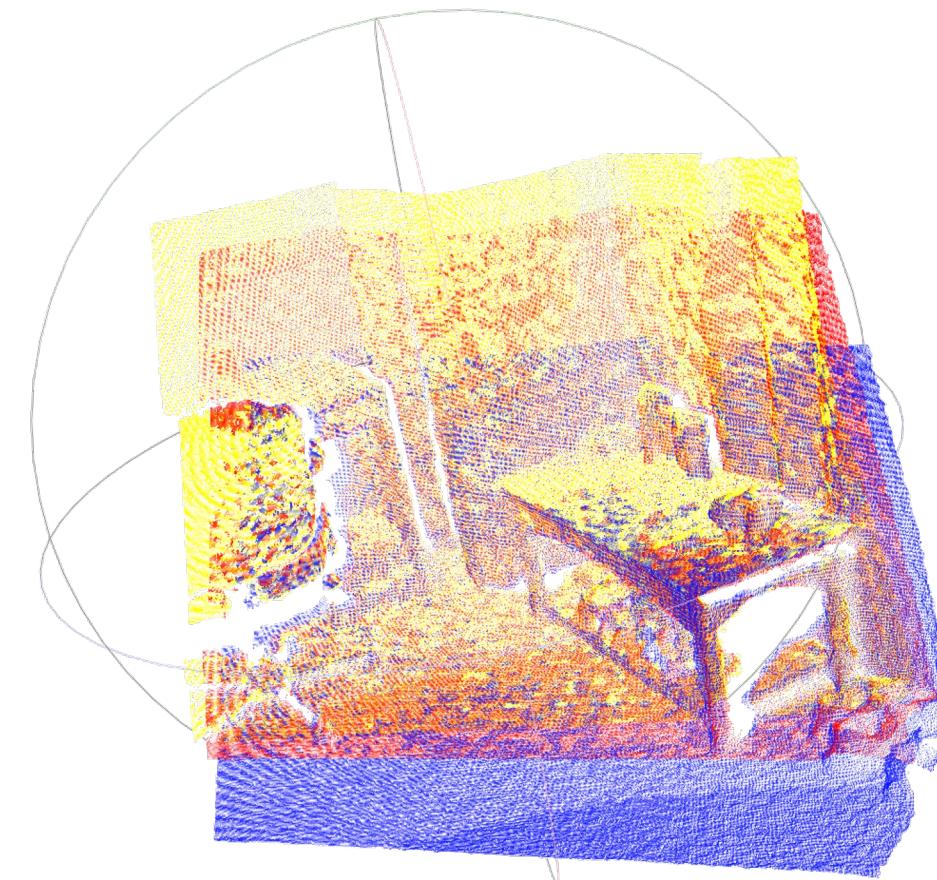
Point Level Contrastive Learning

$$\mathcal{L}_c = - \sum_{(i,j) \in \mathcal{P}} \log \frac{\exp(\mathbf{f}_i \cdot \mathbf{f}_j / \tau)}{\sum_{(\cdot,k) \in \mathcal{P}} \exp(\mathbf{f}_i \cdot \mathbf{f}_k / \tau)}$$

Partial frames from raw RGB-D scan



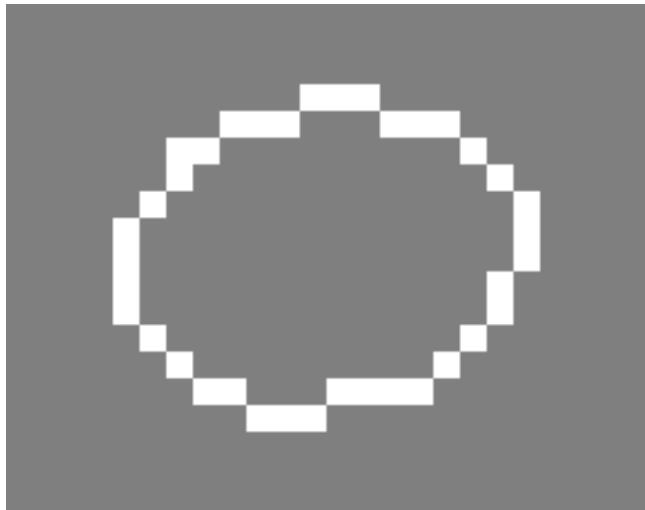
Alignment and correspondences through reconstruction



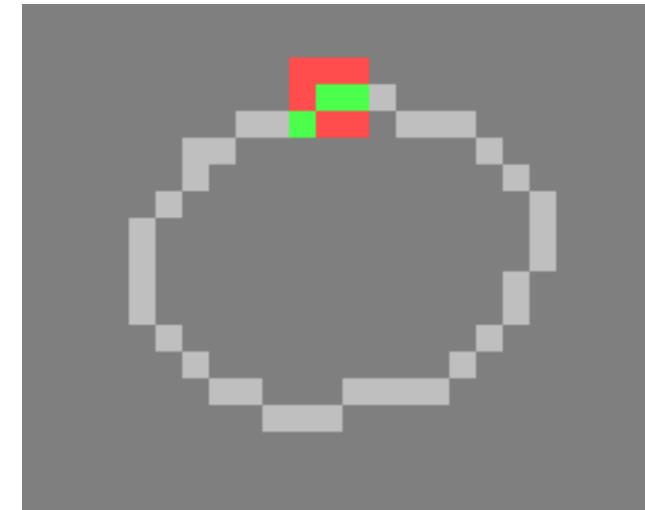
No manual annotation required

Sparse ConvNet (recall: the same thing used in ConvNeXt v2)

Non-zero active sites grow rapidly



Computation only on a sparse set of “active sites”

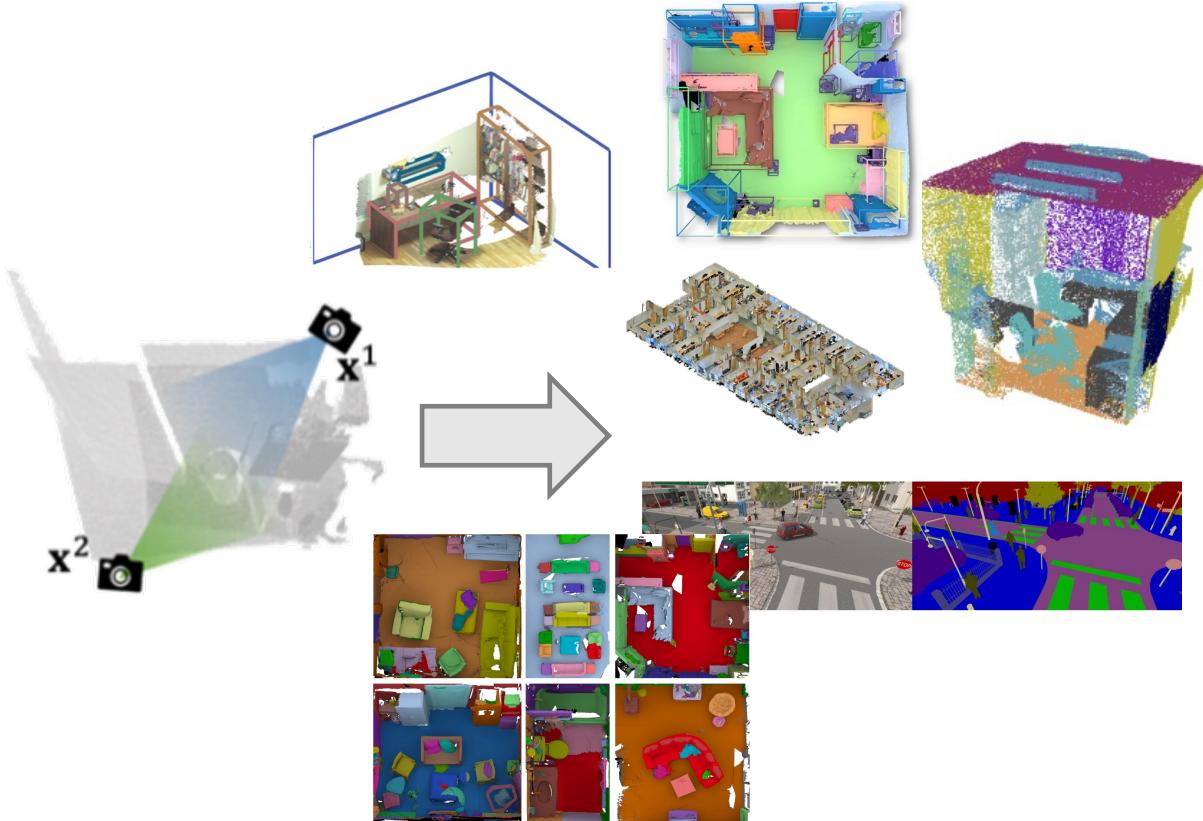


Graham, Benjamin, and Laurens van der Maaten. “Submanifold sparse convolutional networks.” *CVPR 2018*

PointContrast

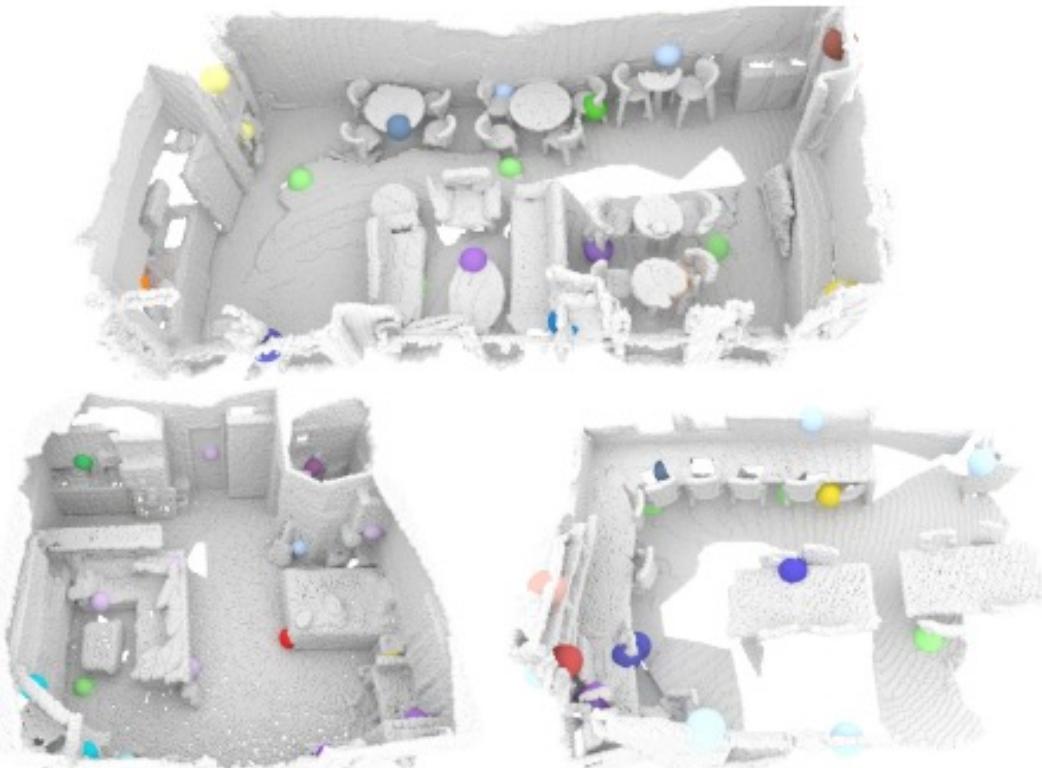
[Xie, Gu, Guo, Qi, Guibas, Litany, ECCV 2020]

First work to demonstrate self-supervised learning is useful for a diverse set of 3D recognition tasks



PointContrast: Downstream Tasks for Fine-tuning					
Datasets	Real / Synth.	Complexity	Env.	Task	Rel. gain
S3DIS	Real	Entire floor, office	Indoor	Segmentation	(+2.7%) mIoU
SUN RGB-D	Real	Medium-sized cluttered rooms	Indoor	Detection	(+3.1%) mAP0.5
ScanNetV2	Real	Large rooms	Indoor	Segmentation	(+1.9%) mIoU
				Detection	(+2.6%) mAP0.5
ShapeNet	Synth.	Single objects	Indoor & outdoor	Classification	(+4.0%) Acc.*
ShapeNetPart	Synth.	Object parts	Indoor & outdoor	Segmentation	(+2.2%) mIoU*
Synthia 4D	Synth.	Street scenes, driving envs.	Outdoor	Segmentation	(+3.3%) mIoU

Representation learning enables data-efficiency



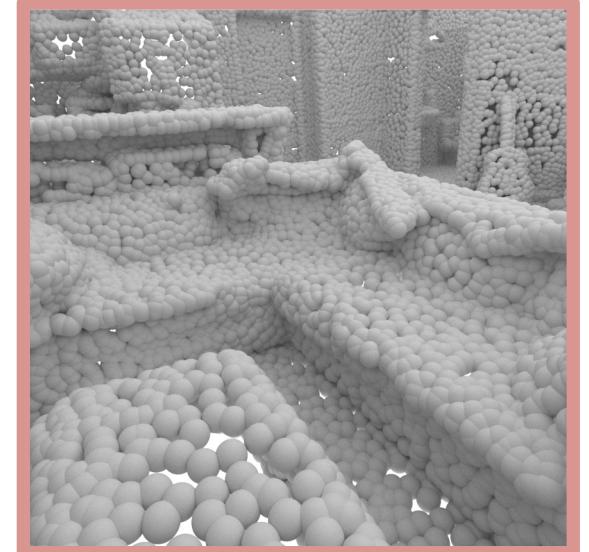
Training Data
with Limited Annotations



Instance Segmentation
Predictions

Using only 20 point labels per scene!

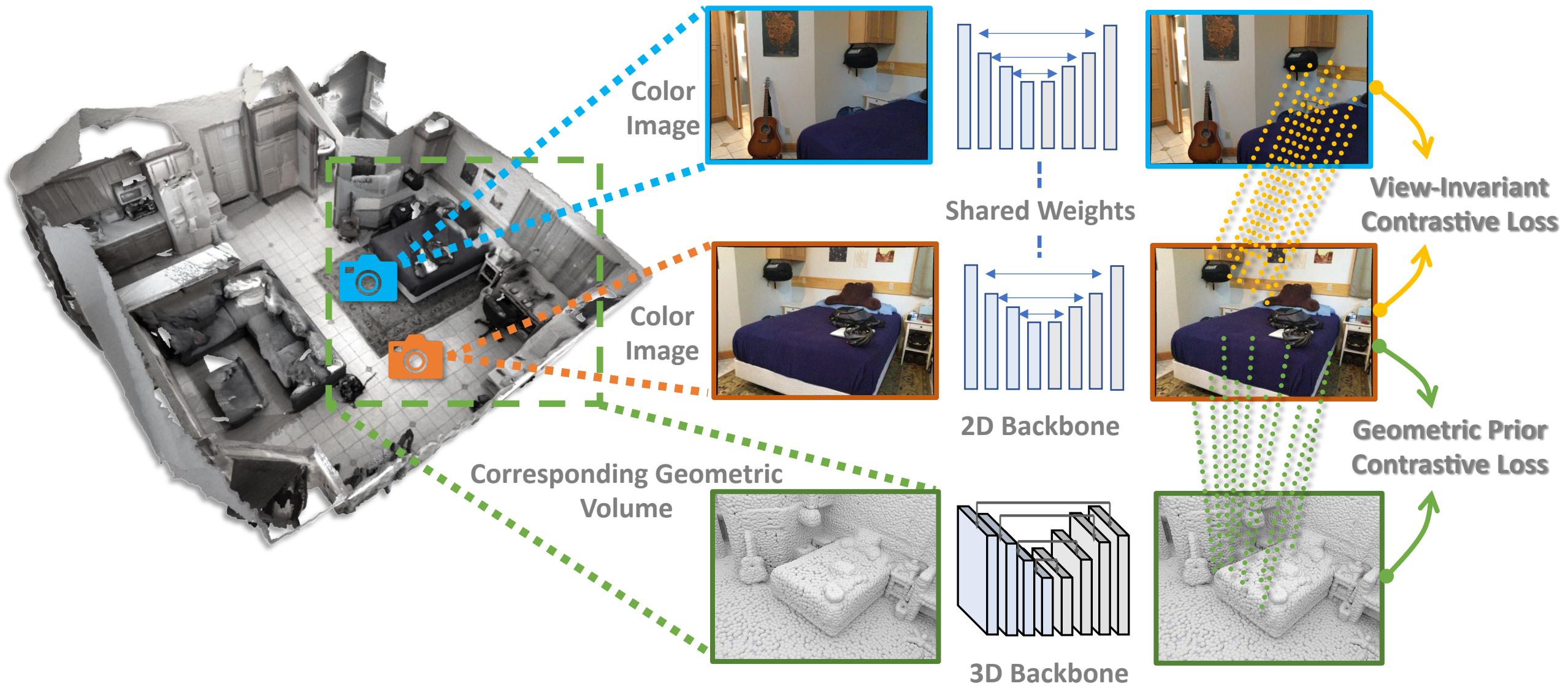
Joint 3D-2D holistic representation learning



- Can 3D prior help 2D tasks?

[Hou, Xie, Graham, Dai, Nießner, ICCV 2021]

Joint View-Invariant and Geometric Prior Contrastive Loss



Connecting representation learning with language



Uncurated Data



Title: Oil spill slick at NSRCC, Tanah Merah with city skyline
Desc: More about the oil spill on the wild shores of singapore blog. 300dpi photo free for download as part of the 2010 International Year of Biodiversity.



Title: Egg tastic!
Desc: we normally write the experation date on our eggs as we do not keep the container. Yesterday I found this little "Easter Egg" left by Erin!



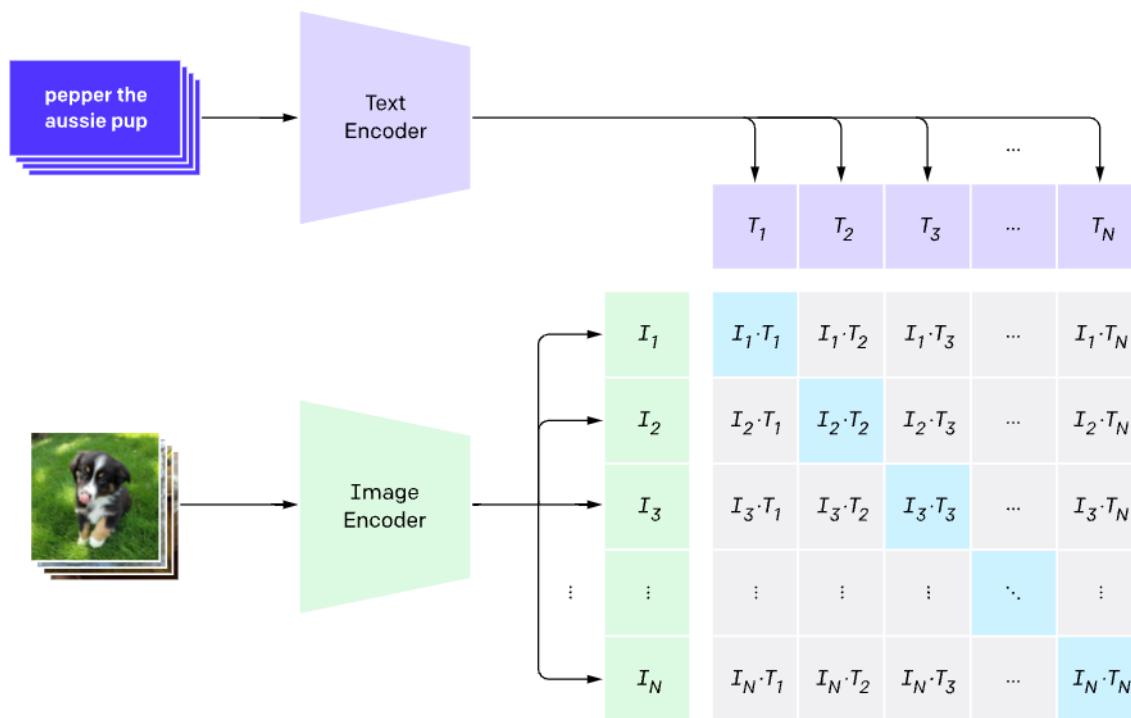
Title: Ping pong project
Desc: Got a set of ping pong balls in for a project. Had to play with them before anything else

Samples from YFCC dataset

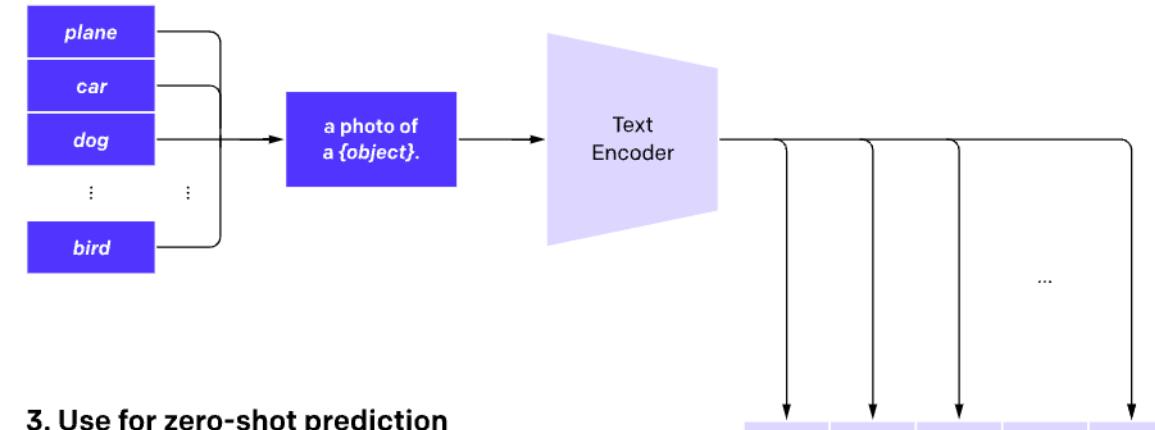
Language: a form of weak supervision?
(Well, not exactly)

CLIP: Connecting text and images

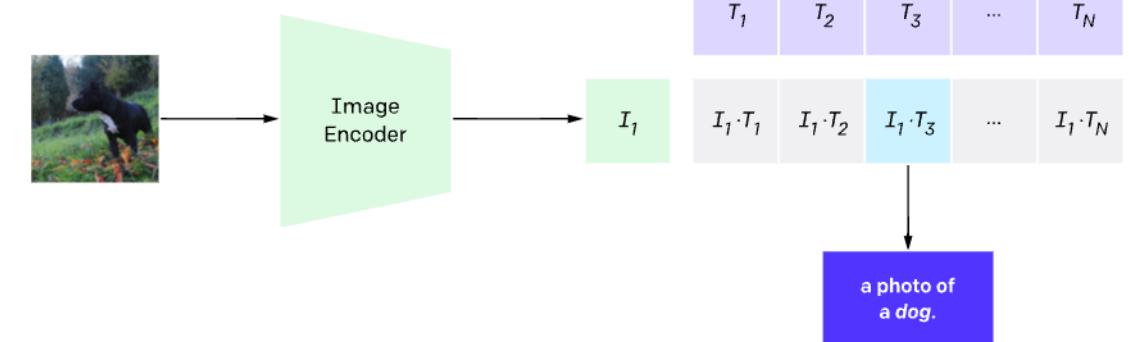
1. Contrastive pre-training



2. Create dataset classifier from label text



3. Use for zero-shot prediction



A data-centric view of model robustness

OpenCLIP

	Dataset Examples			ImageNet	Zero-Shot	ResNet101	CLIP	Δ Score
	ImageNet	ImageNet-V2	ImageNet-R	ObjectNet	ImageNet-Sketch	ImageNet-A	CLIP	
ImageNet								76.2 76.2 0%
ImageNetV2								64.3 70.1 +5.8%
ImageNet-R								37.7 88.9 +51.2%
ObjectNet								32.6 72.3 +39.7%
ImageNet Sketch								25.2 60.2 +35.0%
ImageNet-A								2.7 77.1 +74.4%

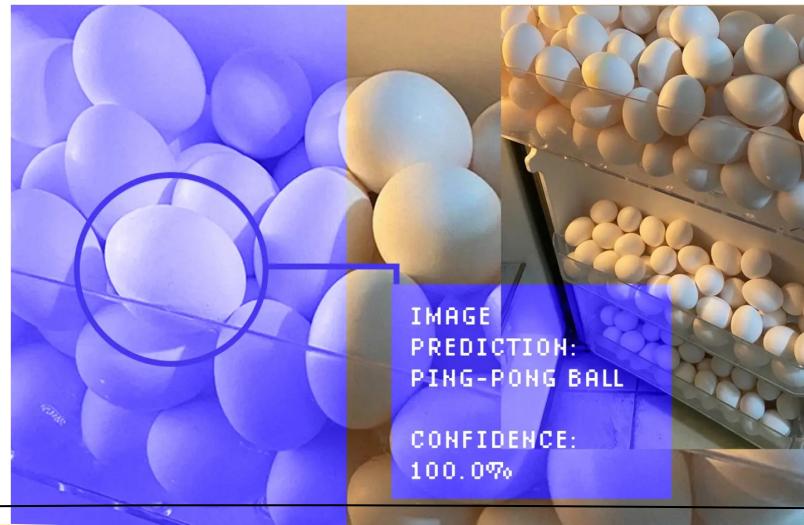
Zero-shot accuracies at resolution 224x224 computed with CLIP Benchmark .		
Dataset	OpenCLIP H/14	OpenCLIP G/14
ImageNet	78.0	80.1
ImageNet-V2	70.8	73.6
ImageNet-R	89.3	92.1
ImageNet-Sketch	66.6	68.9
ObjectNet	69.7	73.0
ImageNet-A	59.2	69.3
CIFAR-10	97.4	98.2
CIFAR-100	84.7	87.5
MNIST	72.9	71.6
SVHN	56.1	62.5
Caltech-101	85.0	86.4
SUN397	75.2	74.5
FGVC Aircraft	42.8	49.7
Country211	30.0	33.8
Cars	93.5	94.6

Building accurate models

AI

The creator of the viral "neural net guesses memes" Twitter account explains how it works

Computer engineer Glen Neumann breaks down ResNeXtGuesser, his bot that often hilariously reacts to bizarre images.



"Just as artificial intelligence notoriously lacks common sense, it also lacks an understanding of memes."



neural net guesses memes
@ResNeXtGuesser

...

Image prediction: ping-pong ball
Confidence: 99.99%
Submission by @Minish900



10:53 PM · Nov 9, 2021 · ImageGuesserBot

10.5K Retweets 1,181 Quote Tweets 115.3K Likes



Recall: model robustness from model scaling & self-supervised learning

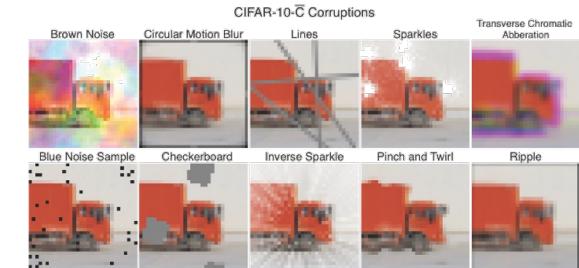
Masked Autoencoder

dataset	ViT-B	ViT-L	ViT-H	ViT-H ₄₄₈	prev best
IN-Corruption ↓ [27]	51.7	41.8	33.8	36.8	42.5 [32]
IN-Adversarial [28]	35.9	57.1	68.2	76.7	35.8 [41]
IN-Rendition [26]	48.3	59.9	64.4	66.5	48.7 [41]
IN-Sketch [60]	34.5	45.3	49.6	50.9	36.0 [41]

Scaling brings +40% in accuracy

ConvNeXt

Model	Data/Size	FLOPs / Params	Clean	C (↓)	Ā (↓)	A	R	SK
ResNet-50	1K/224 ²	4.1 / 25.6	76.1	76.7	57.7	0.0	36.1	24.1
Swin-T [42]	1K/224 ²	4.5 / 28.3	81.2	62.0	-	21.6	41.3	29.1
RVT-S* [44]	1K/224 ²	4.7 / 23.3	81.9	49.4	37.5	25.7	47.7	34.7
ConvNeXt-T	1K/224 ²	4.5 / 28.6	82.1	53.2	40.0	24.2	47.2	33.8
Swin-B [42]	1K/224 ²	15.4 / 87.8	83.4	54.4	-	35.8	46.6	32.4
RVT-B* [44]	1K/224 ²	17.7 / 91.8	82.6	46.8	30.8	28.5	48.7	36.0
ConvNeXt-B	1K/224 ²	15.4 / 88.6	83.8	46.8	34.4	36.7	51.3	38.2
ConvNeXt-B	22K/384 ²	45.1 / 88.6	86.8	43.1	30.7	62.3	64.9	51.6
ConvNeXt-L	22K/384 ²	101.0 / 197.8	87.5	40.2	29.9	65.5	66.7	52.8
ConvNeXt-XL	22K/384 ²	179.0 / 350.2	87.8	38.8	27.1	69.3	68.2	55.0



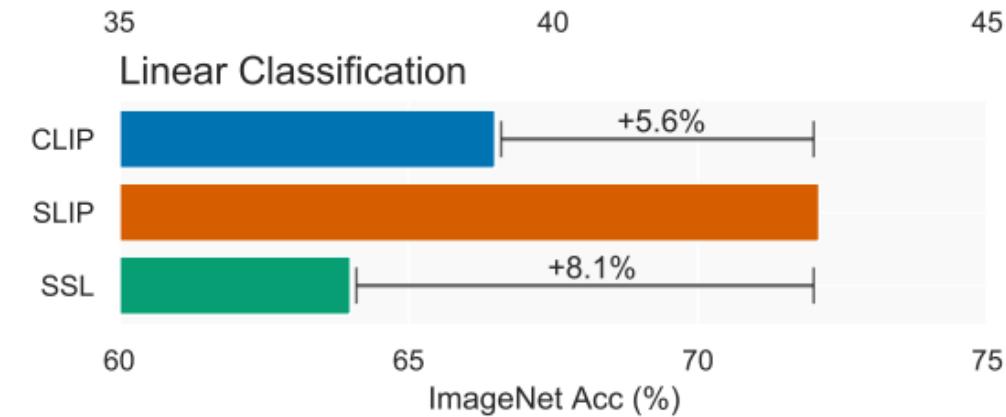
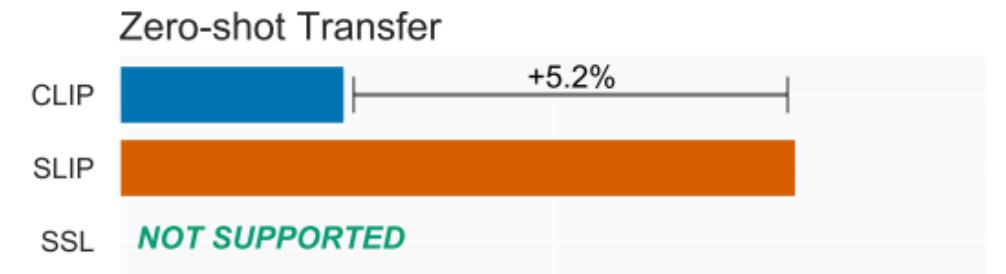
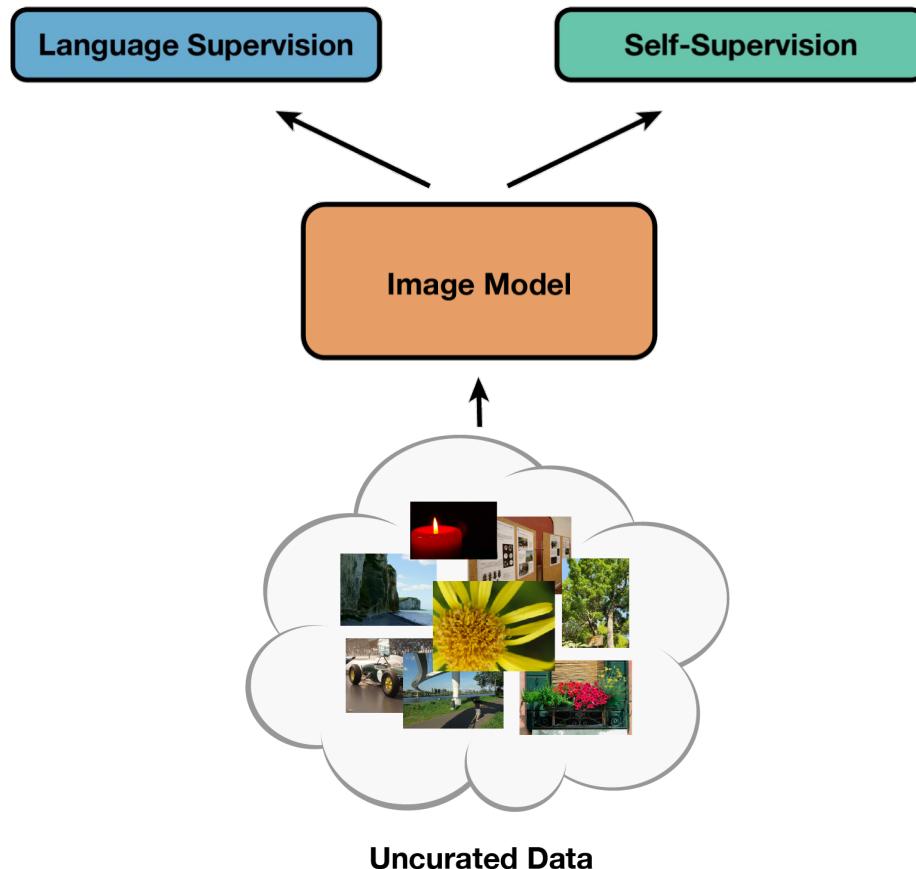
[Mintun, Kirillove, Xie, NeurIPS 2021]

ConvNeXt – similar story

These are different concepts and should be studied independently.

Language-supervision meets self-supervision

[Mu, Kirillov, Wagner, Xie, arXiv 2021]



They are complementary and
we can have the best of both worlds!

Curation-in-Training: Co-designing data and objective

[Hu, Xie, et al. ArXiv 2022]

Data quality matters

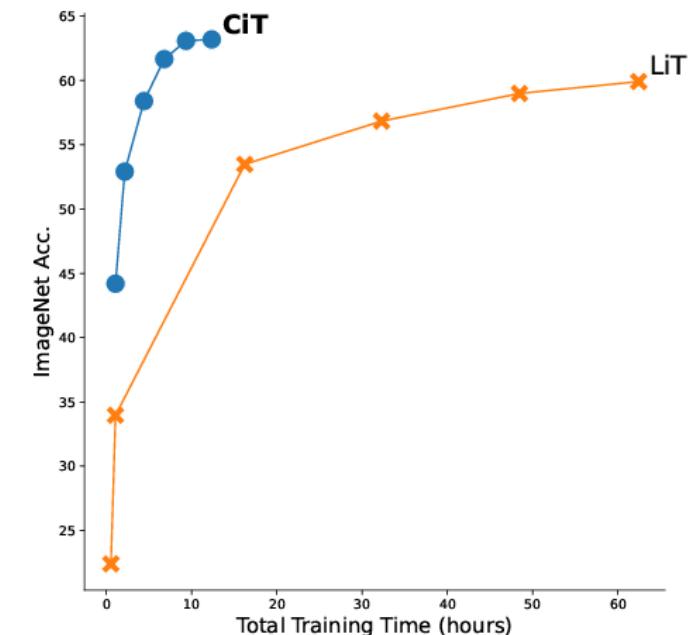
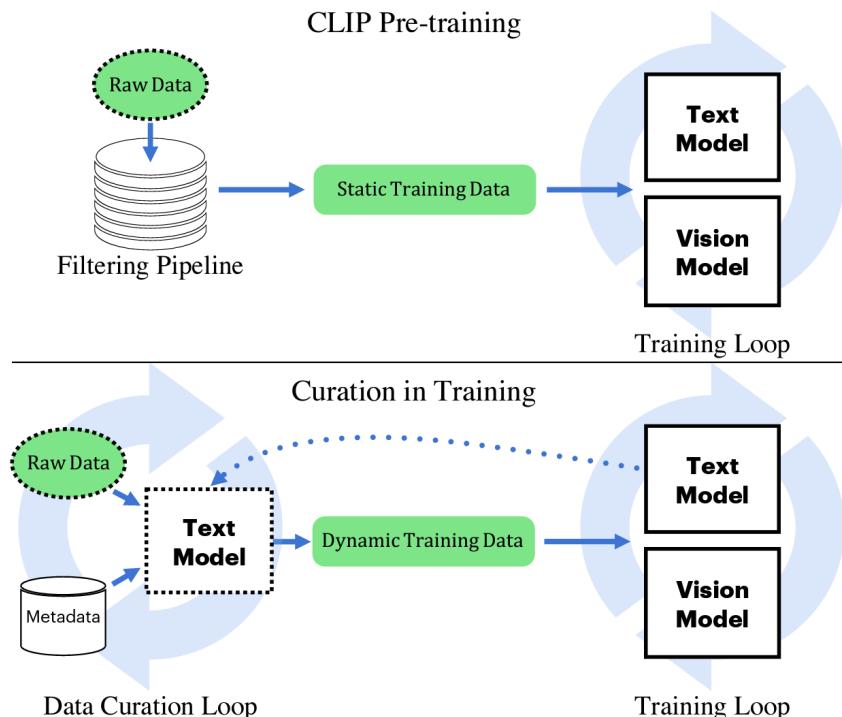
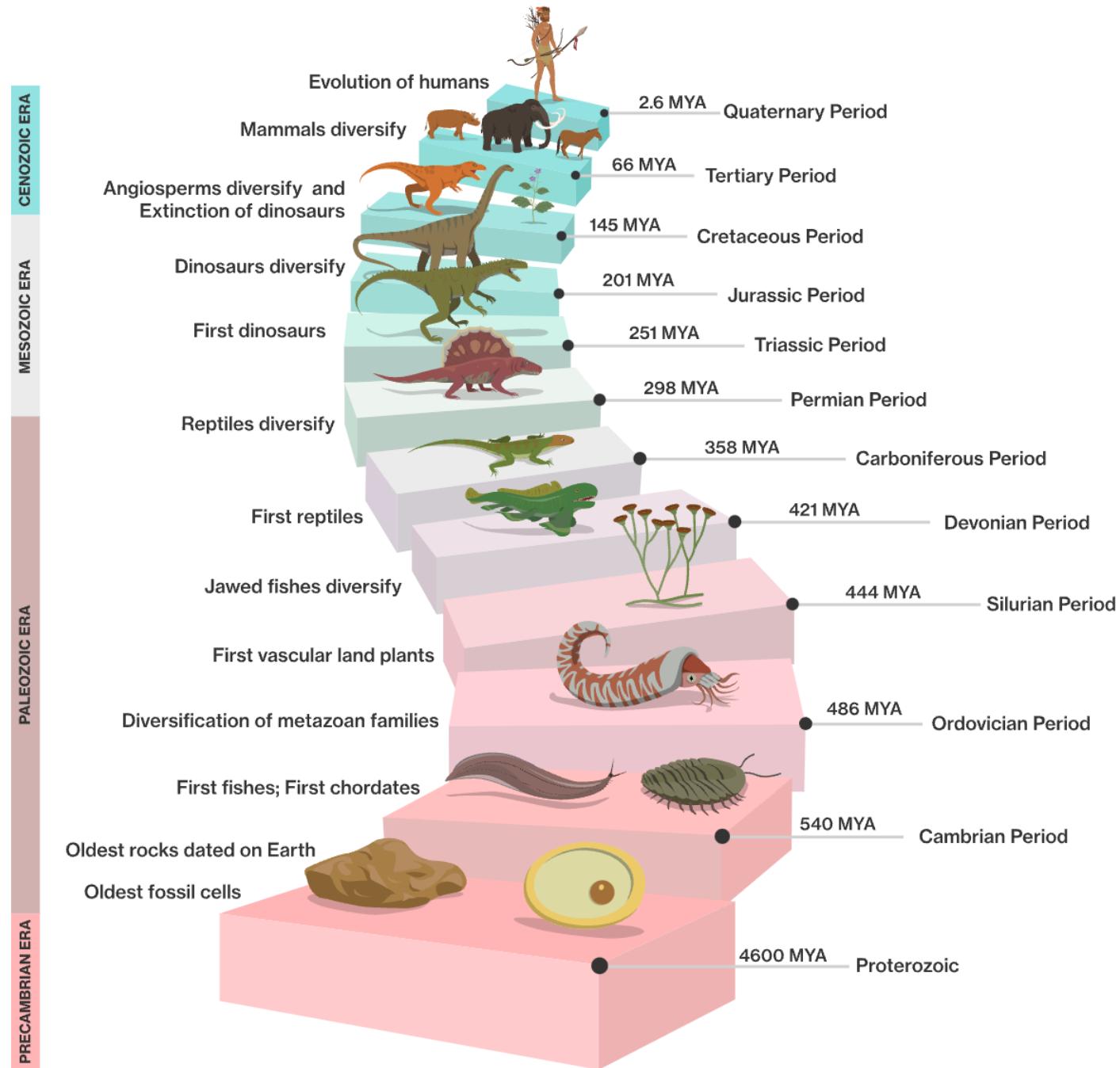


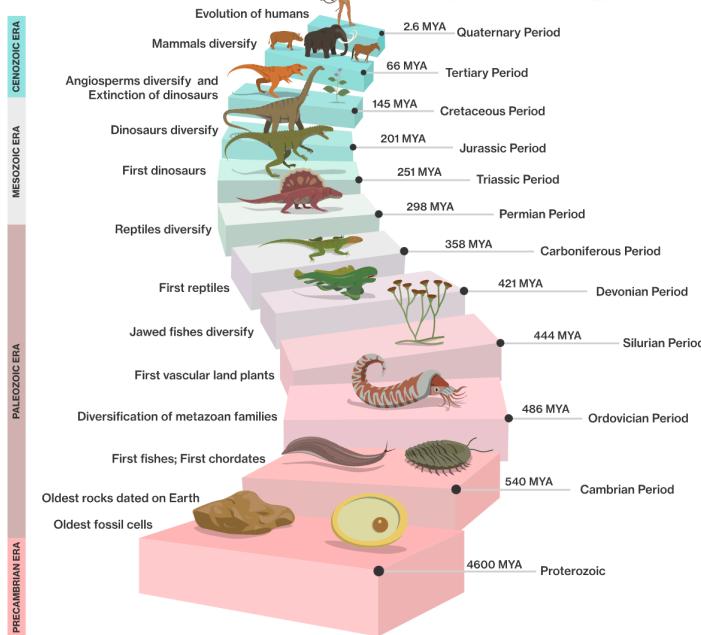
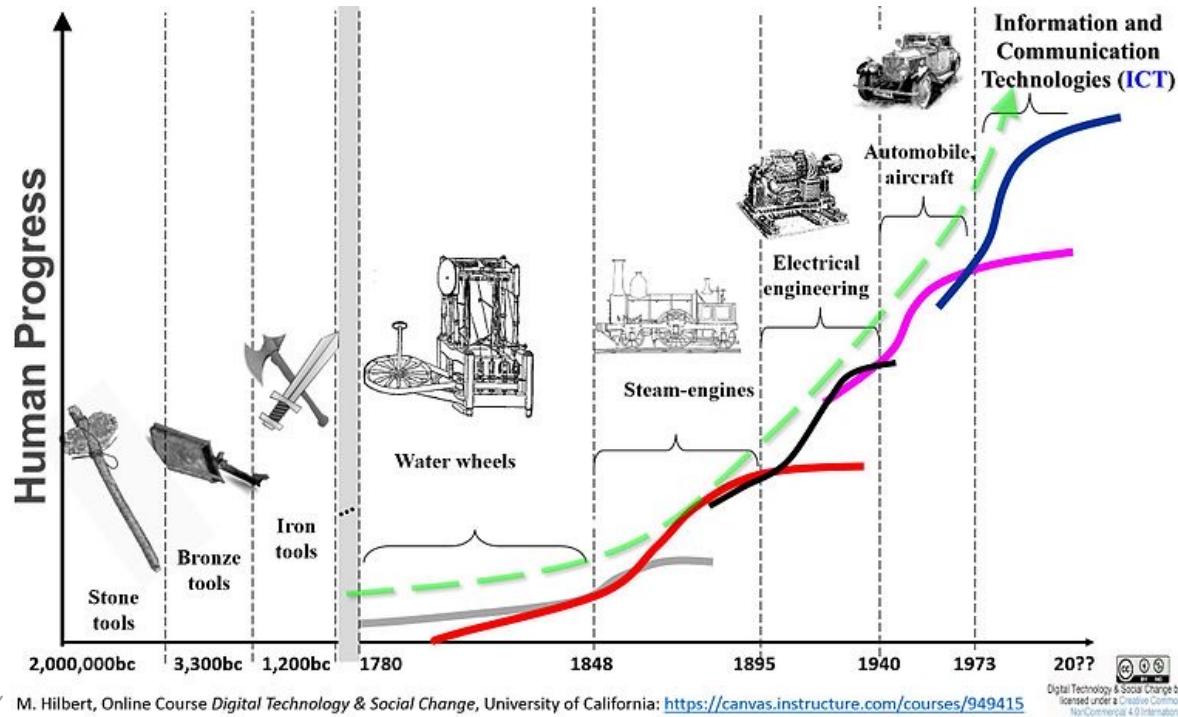
Figure 2. CiT on provides $>5\times$ speedup and +3.4% accuracy gain over LiT [38] on AugReg ViT-B/32 vision encoders. Training data is YFCC15M. Models are evaluated at 6 evenly sampled iterations.

Architecture / **O**bjective / **D**ata

To summarize:

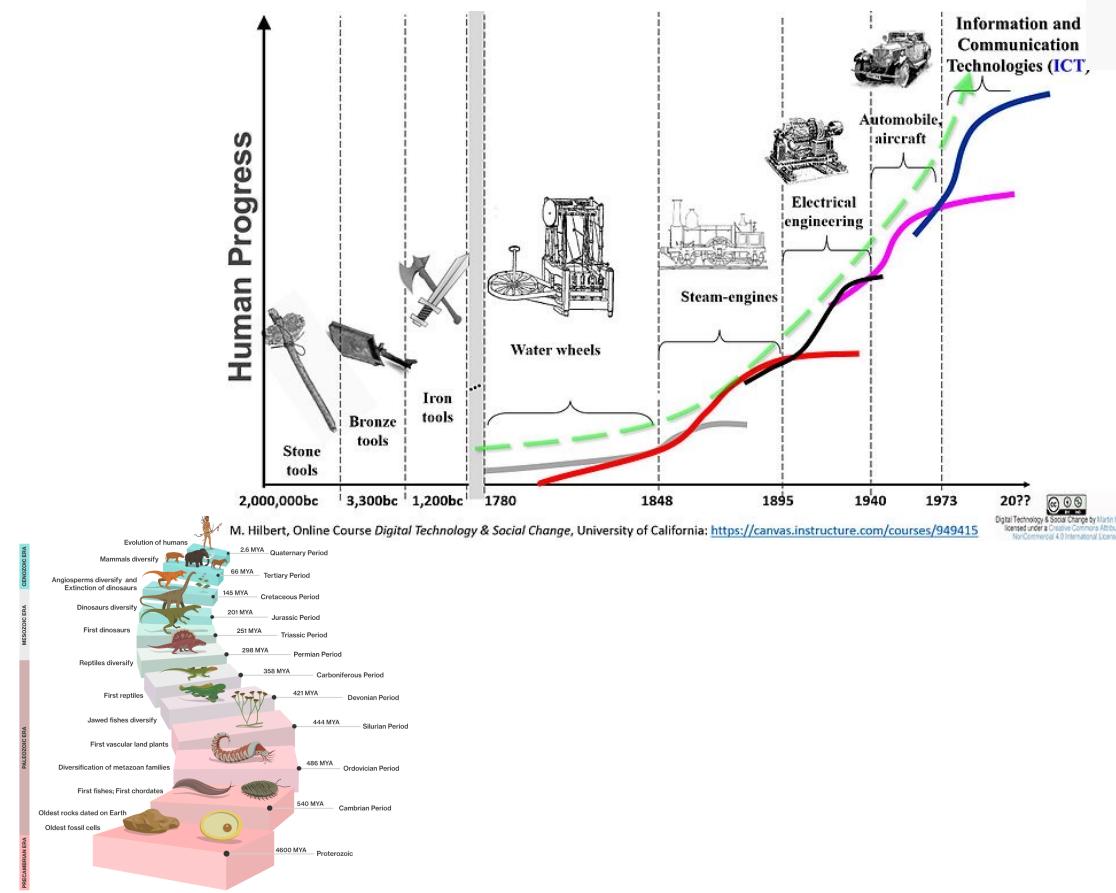
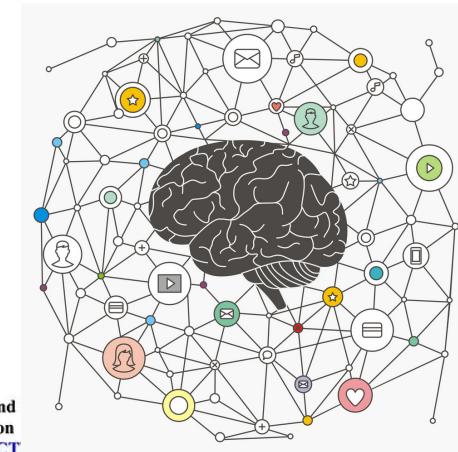
- Co-design matters.
- Implication to model robustness.
- Diverse tasks, diverse applications.
- Great progress, but do we have vision foundation models?





Digital Technology & Social Change by Martin Hilbert
licensed under a Creative Commons Attribution-NonCommercial 4.0 International License.

Internet





Geoffrey Hinton
@geoffreyhinton

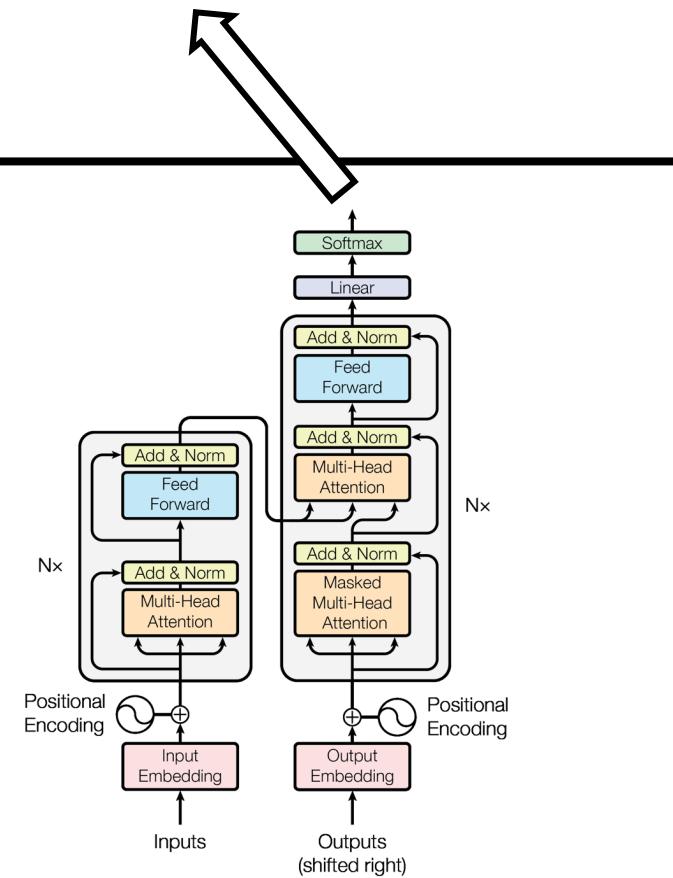
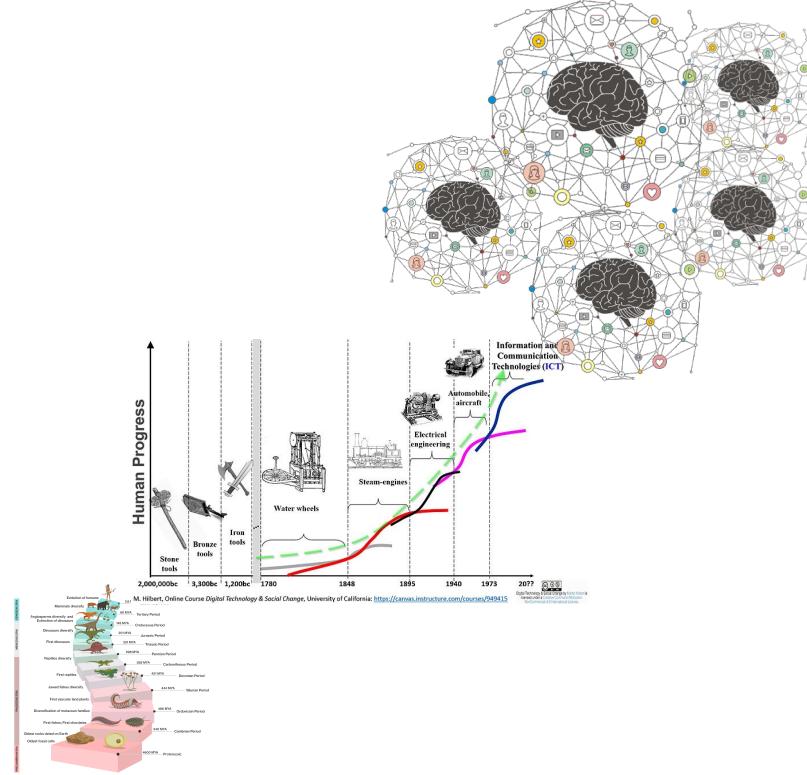
Caterpillars extract nutrients which are then converted into butterflies.
People have extracted billions of nuggets of understanding and GPT-4 is
humanity's butterfly.

4:27 PM · Mar 14, 2023 · 452.2K Views

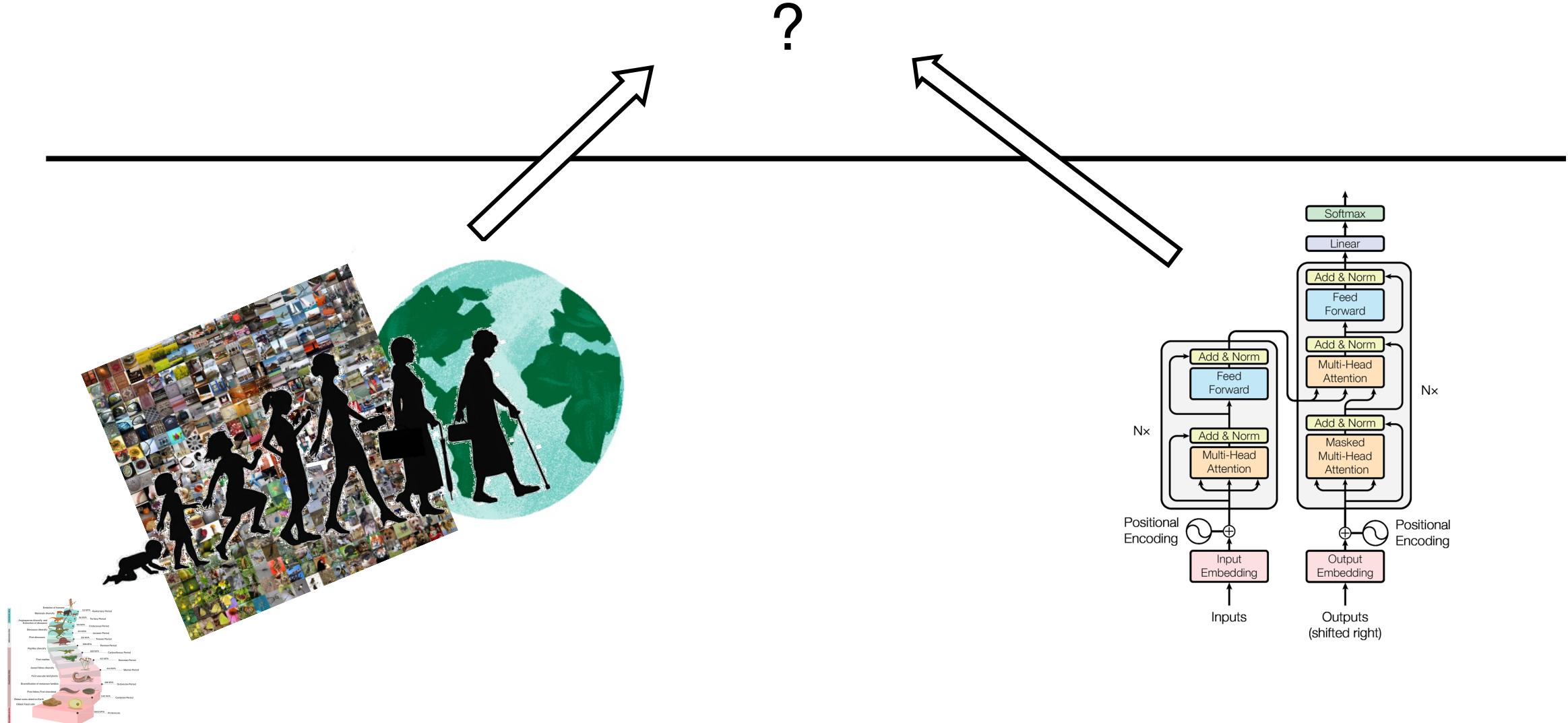
453 Retweets 103 Quotes 2,435 Likes 191 Bookmarks

Next token prediction

$$L_1(\mathcal{U}) = \sum_i \log P(u_i | u_{i-k}, \dots, u_{i-1}; \Theta)$$



For *robust* vision: shortcut no more...



Architecture / **O**bjective / **D**ata

Thank You!