

Convergence of Architectures and Learning Methods in Computer Vision with NLP and Other AI fields

Han Hu

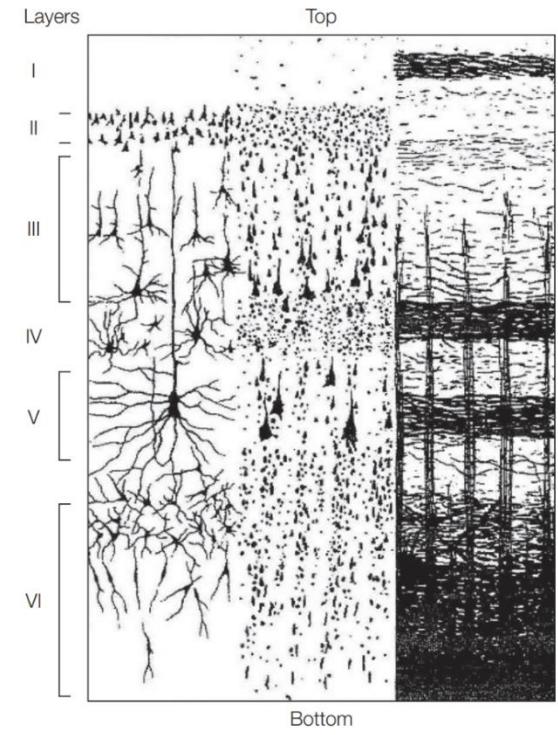
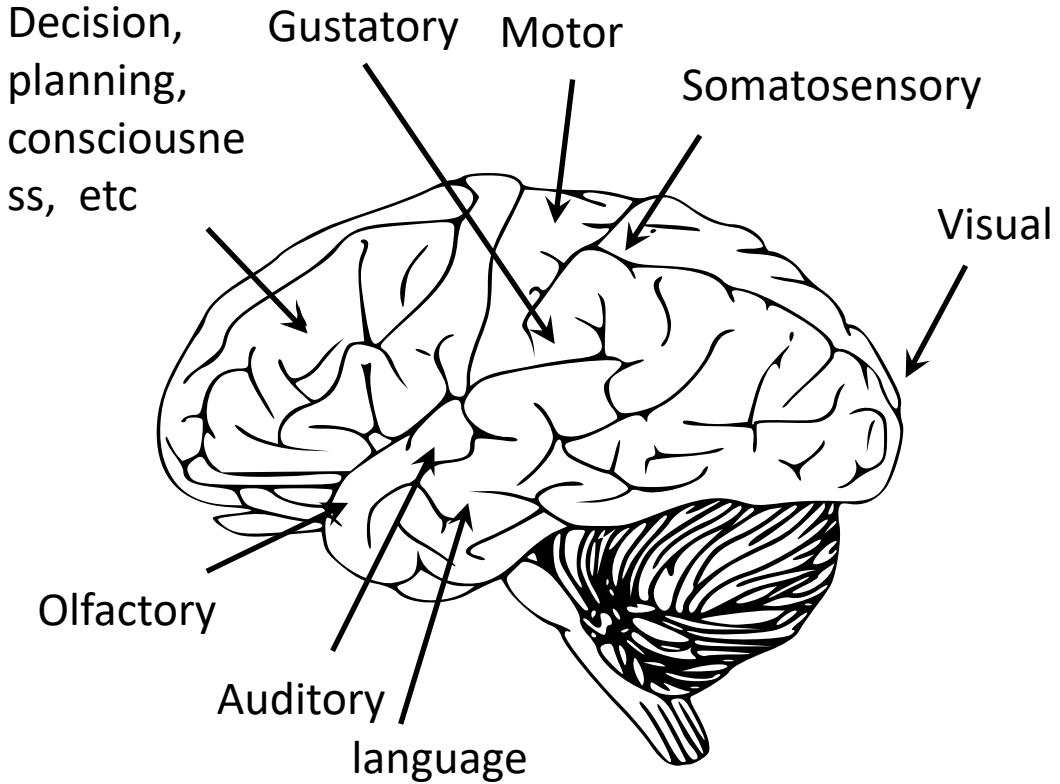
Microsoft Research Asia

CVPR2022 OpenMMLab Tutorial

June 20th , 2022

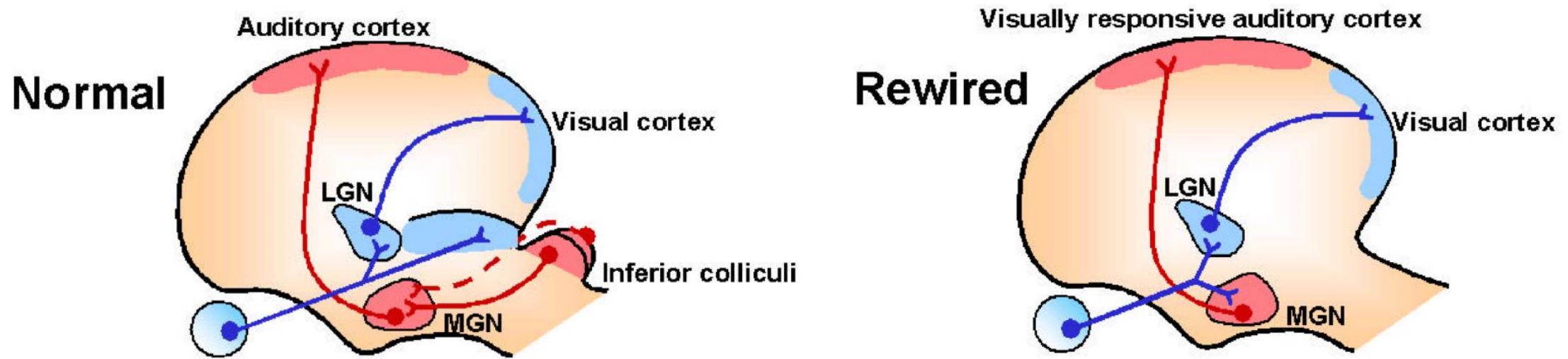
Human brain: A **universal architecture** and universal learning machine

- Universal architecture (innate): neocortex



Human brain: A **universal architecture** and universal learning machine

- Plasticity of neocortex: a mouse brain experiment



Human brain: A universal architecture and universal learning machine

- Learning (nurture): Yann LeCun's cake

- ▶ “Pure” Reinforcement Learning (**cherry**)
 - ▶ The machine predicts a scalar reward given once in a while.
 - ▶ **A few bits for some samples**
- ▶ Supervised Learning (**icing**)
 - ▶ The machine predicts a category or a few numbers for each input
 - ▶ Predicting human-supplied data
 - ▶ **10→10,000 bits per sample**
- ▶ Self-Supervised Learning (**cake génoise**)
 - ▶ The machine predicts any part of its input for any observed part.
 - ▶ Predicts future frames in videos
 - ▶ **Millions of bits per sample**



Credit: Yann LeCun

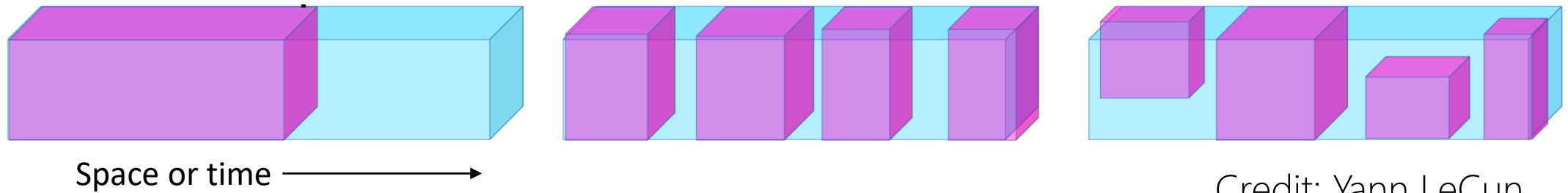
Human brain: A universal architecture and **universal learning** machine

- Learning (nurture) : baby's learning through observation (self-supervised)



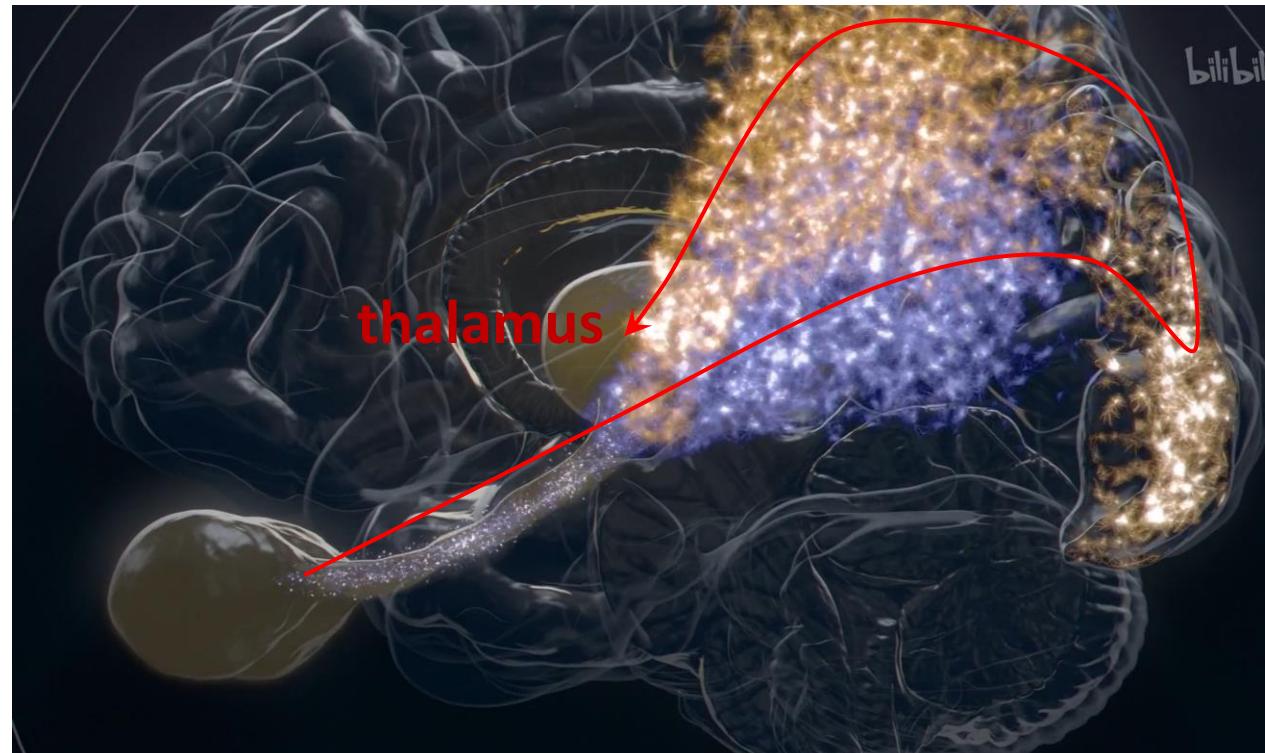
Human brain: A universal architecture and **universal learning** machine

- Unified learning approach: Compare difference between **prediction** and ground truth input
 - Vision, touch, sound, language, etc



Human brain: A universal architecture and universal learning machine

- Unified learning approach: Compare difference between **prediction** and ground truth input
 - Thalamus plays a key role

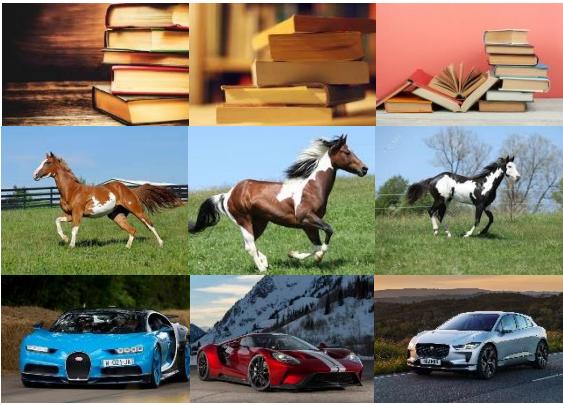


Credit: David Eagleman

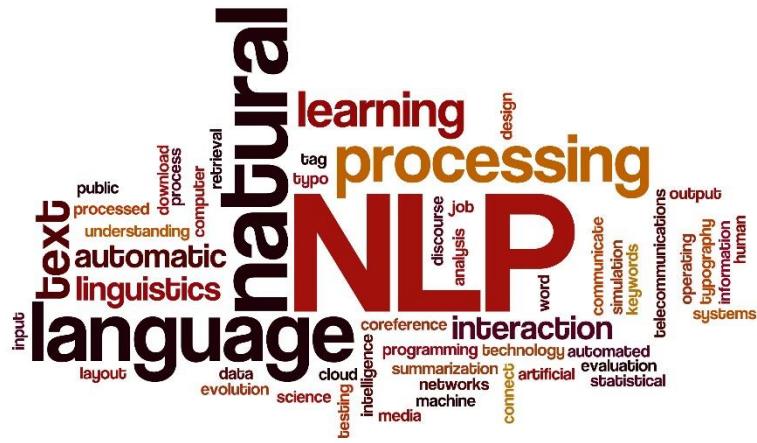
Towards Converged AI Architectures

(Swin) Transformer

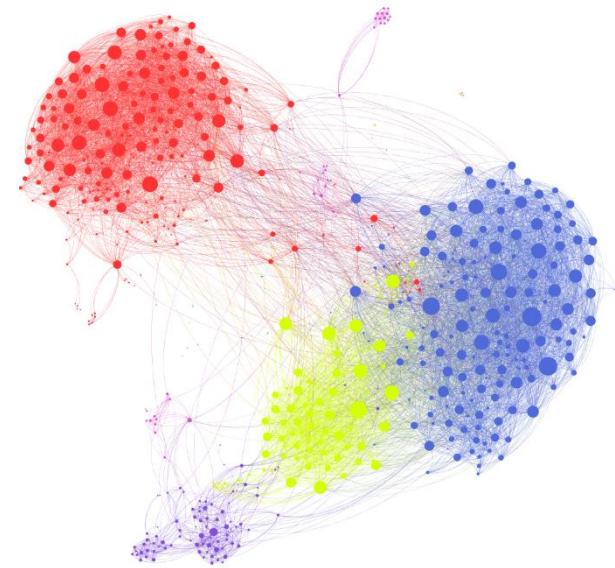
Mainstream models of AI sub-fields



vision-CNN



language-Transformer



social-Graph Networks

Model evolution for NLP or sequential data

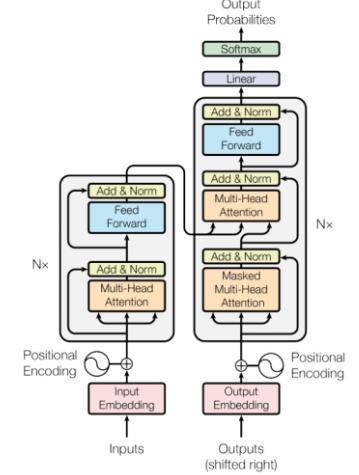
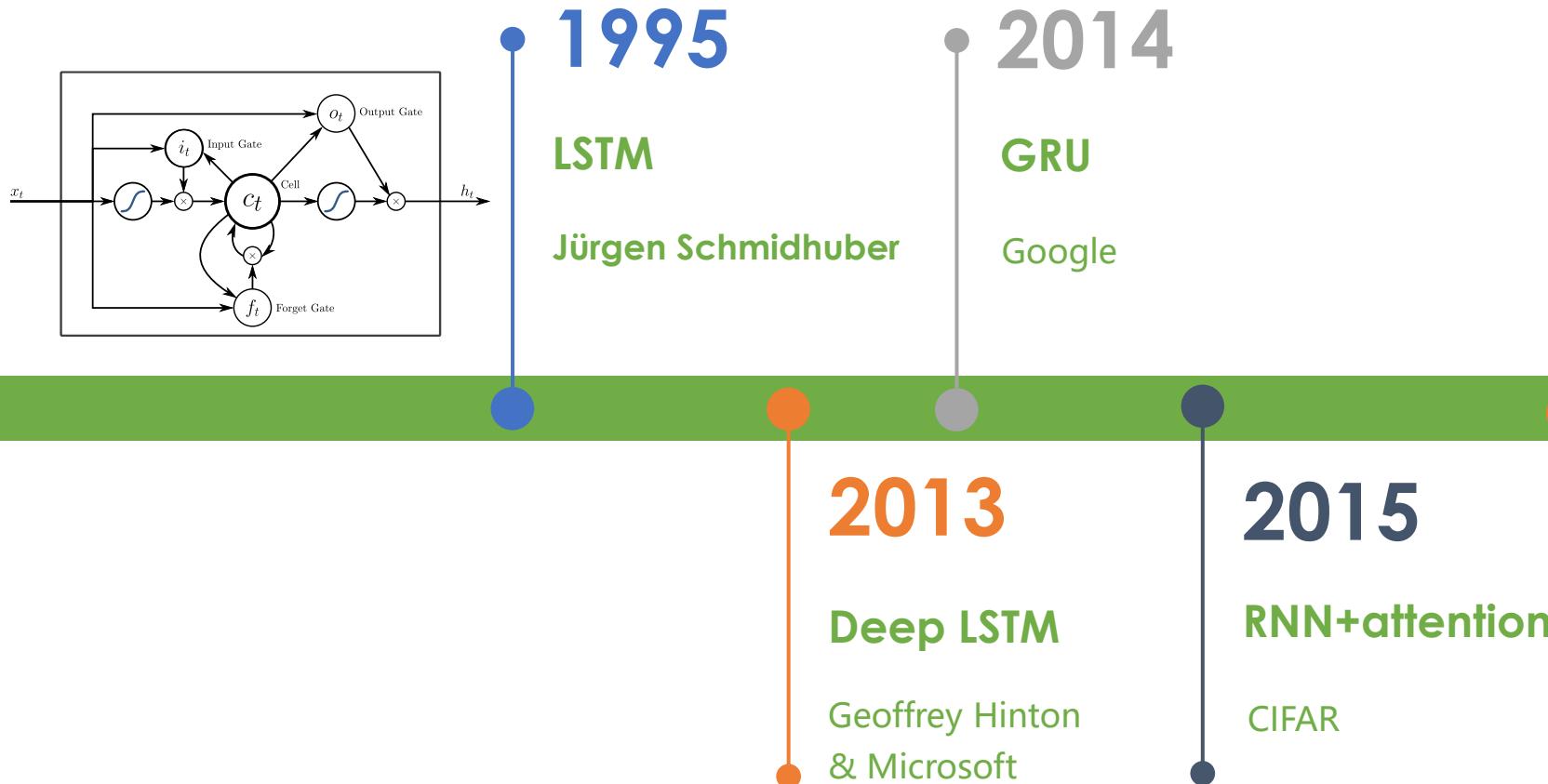


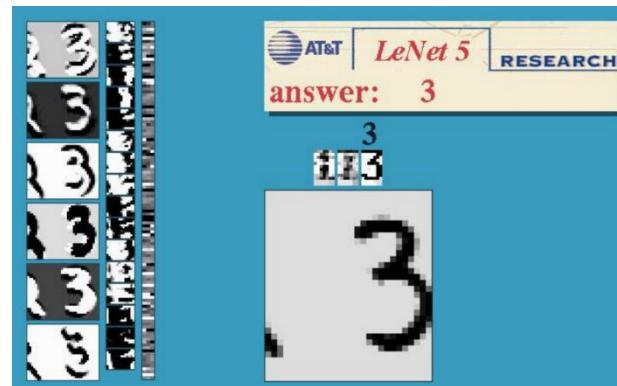
Figure 1: The Transformer - model architecture.

Model evolution for CV

• 1980

CNN

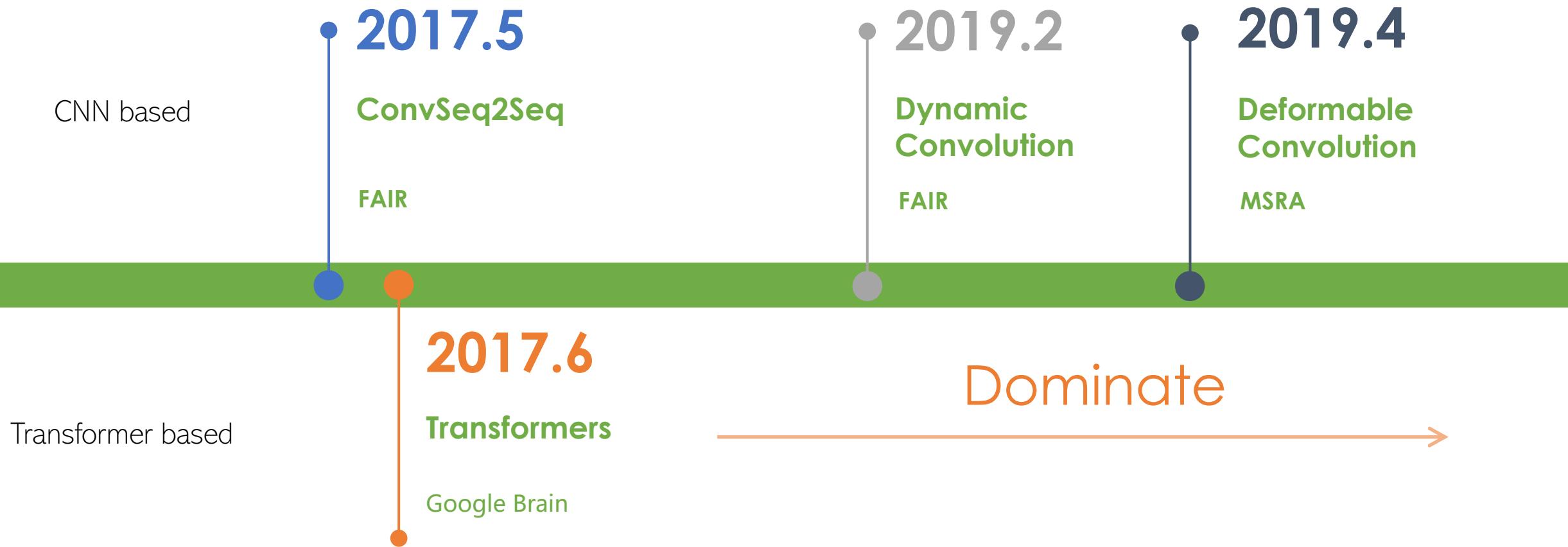
Kunihiko Fukushima
Yann LeCun



AlexNet, GoogleNet, VGGNet, ResNet, DenseNet, ...

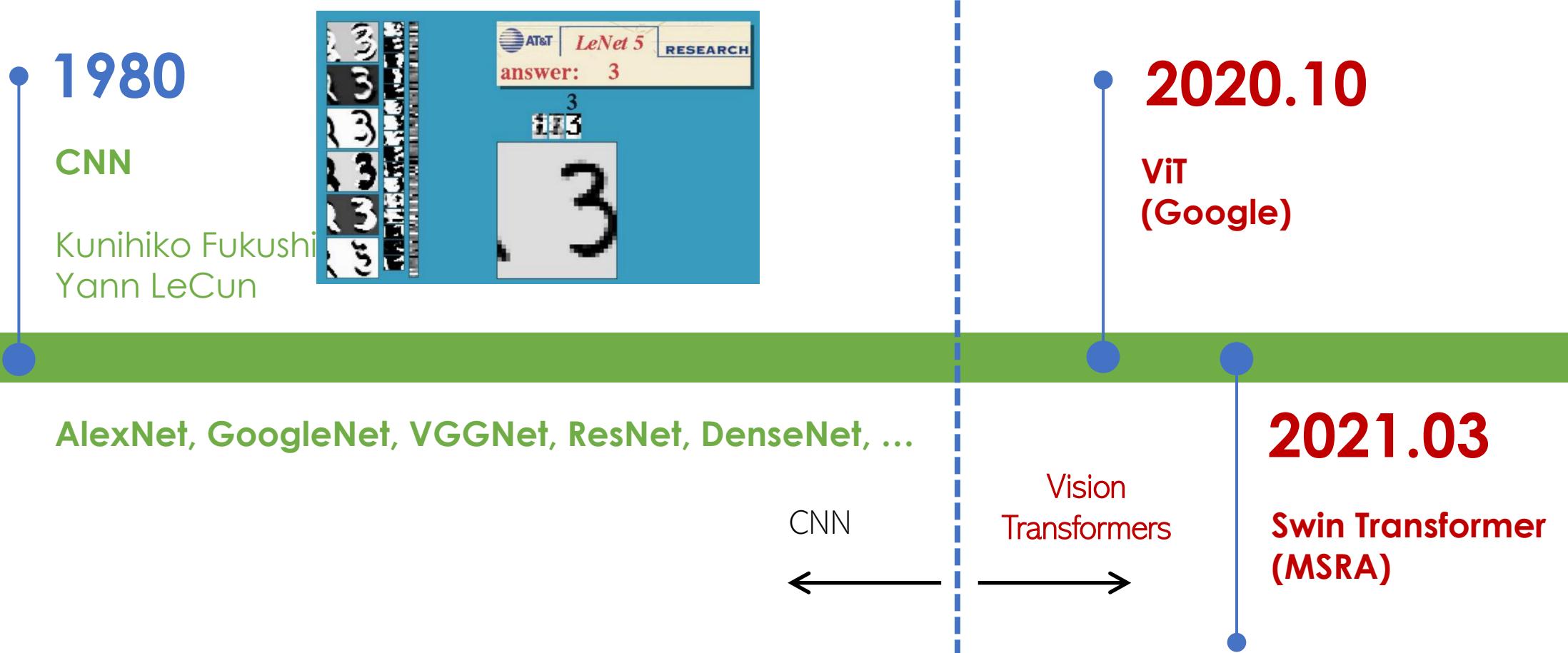
Can NLP and CV use the same architectures?

- Adapt CNN or convolution for NLP



Can NLP and CV use the same architectures?

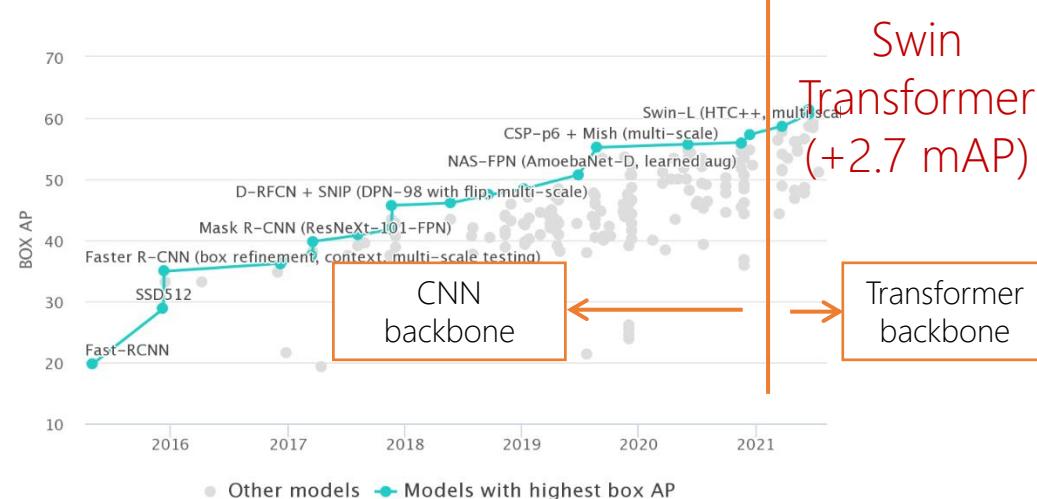
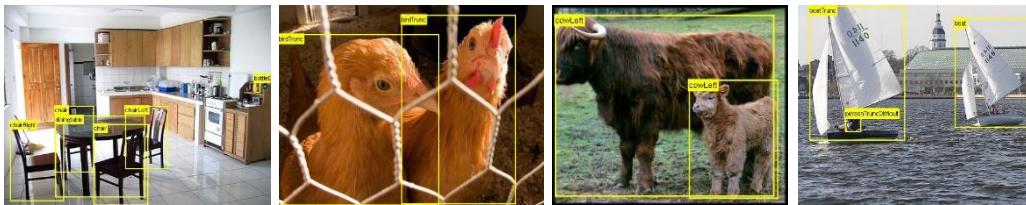
- Adapt Transformers or attention for CV



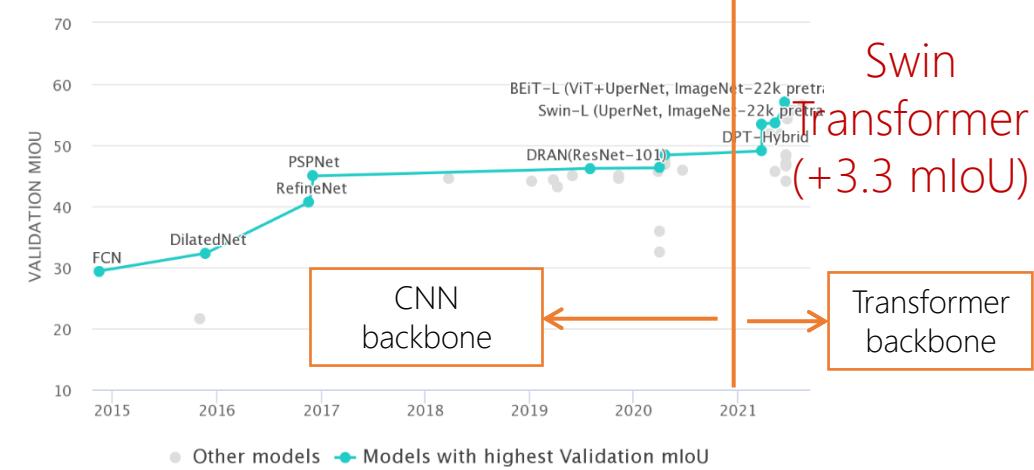
Swin Transformer: A general-purpose vision backbone

- The first time to significantly surpass previous best CNN on broad vision tasks

COCO object detection

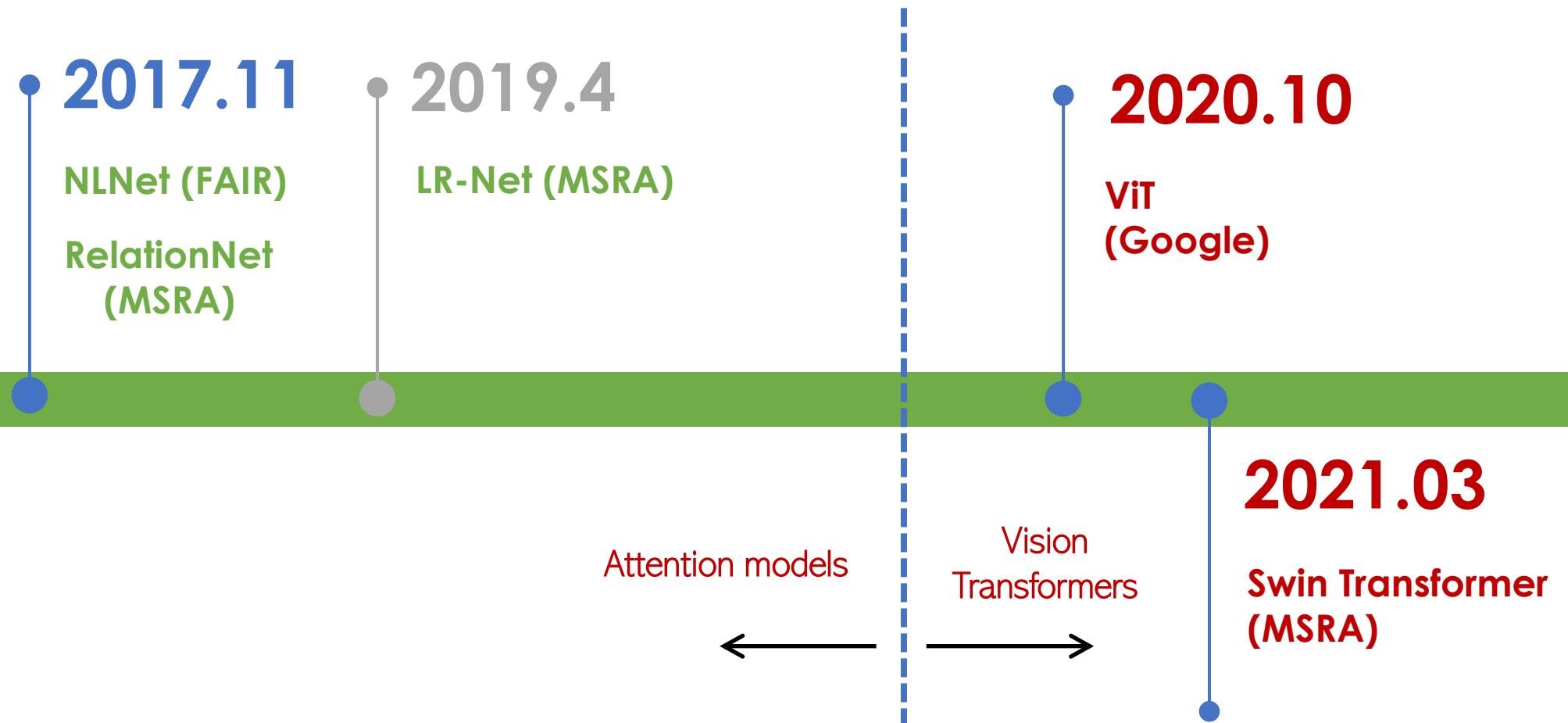


ADE20K semantic segmentation



Can NLP and CV use same architectures?

- Vision Transformers and attention models



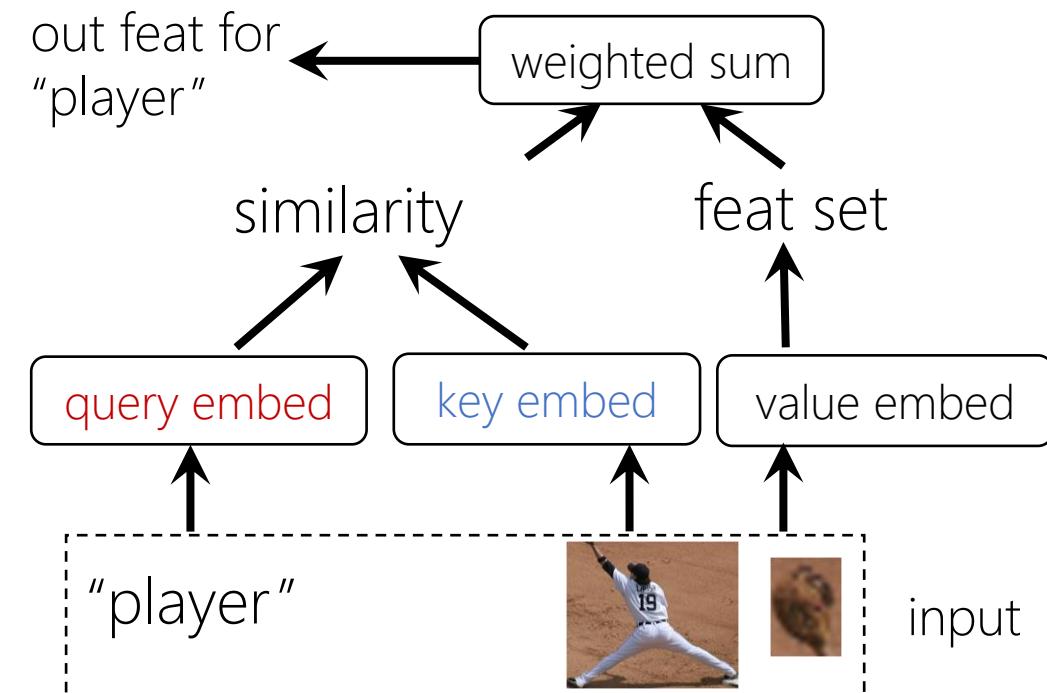
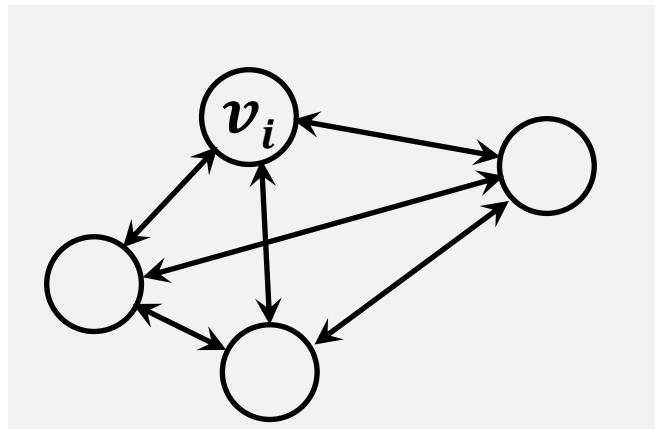
It is a long-term belief

Towards Universal Learning
Machine: Self-Attention for Visual
Modeling

2019.7 @valse webinar

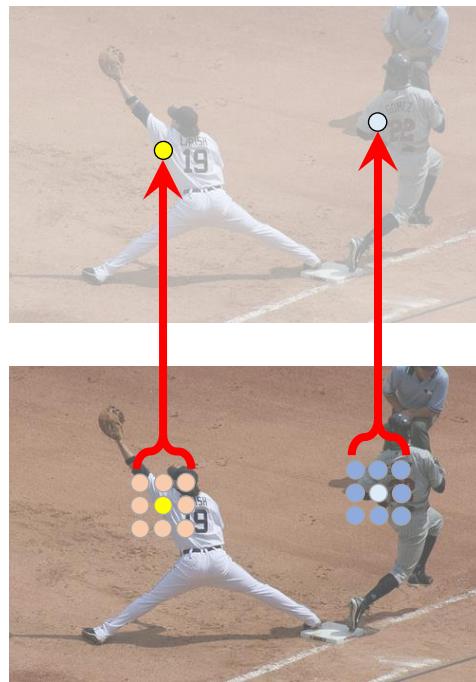
Why is it Transformer or attention?

- General modeling capability
 - All concepts (concrete or abstract) and their relationships can be modeled by a graph
 - Modeling arbitrary relationship via verification, which is hard for CNN

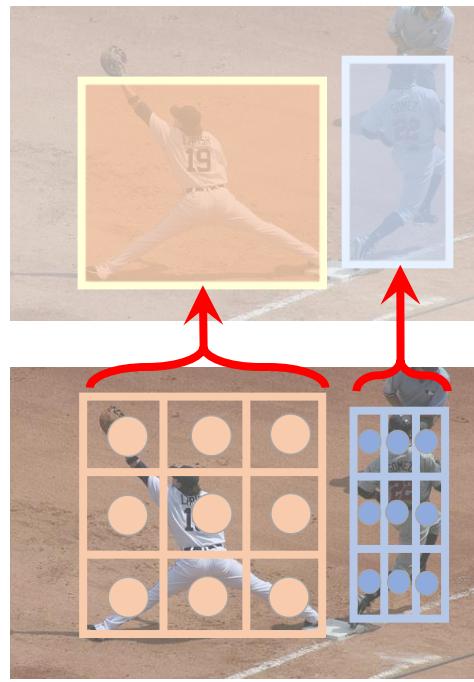


Transformer is so general

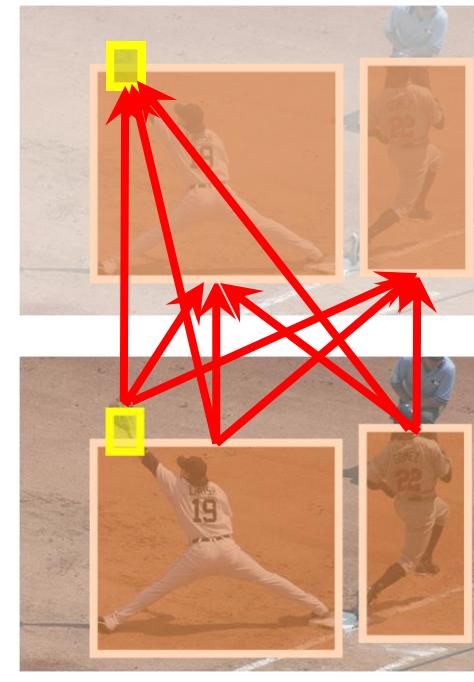
- Can model all of pixel-to-pixel, object-to-pixel, object-to-object relationships, which covers almost all vision applications



pixel-to-pixel



object-to-pixel



object-to-object

Our works (2017-2021)

- Pixel-to-pixel

- [1] GCNet: Non-local networks meet squeeze-excitation networks and beyond, ICCVW 2019
- [2] Local Relation Networks for Image Recognition. ICCV 2019 (**the first full-attention visual backbone**)
- [3] Disentangled Non-local Neural Networks. ECCV 2020
- [4] Swin Transformer. ICCV 2021 (**best paper award**)
- [5] Swin Transformer V2. CVPR 2022

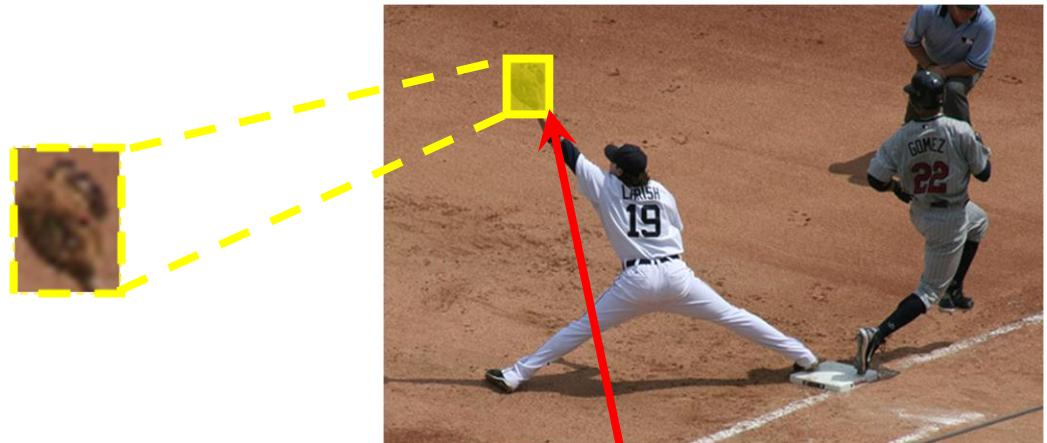
- Object-to-pixel

- [6] Learning Region Features for Object Detection. ECCV 2018
- [7] RelationNet++: Bridging Visual Representations for Object Detection via Transformer Decoder, NeurIPS 2020

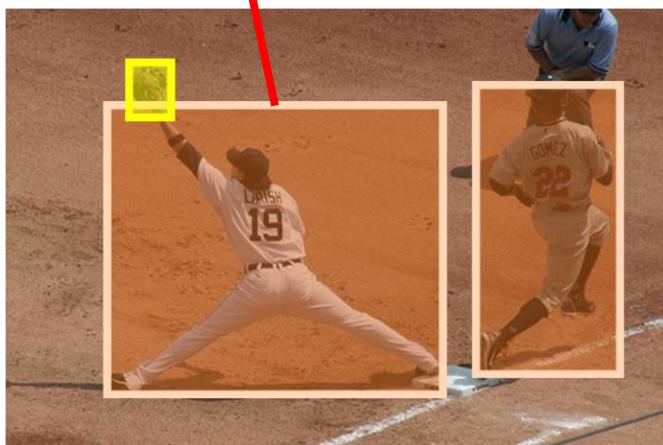
- Object-to-object

- [8] Relation Networks for Object Detection, CVPR2018 (**the first end-to-end object detector**)
- [9] Spatial-temporal relation networks for multi-object tracking. ICCV 2019
- [10] Memory Enhanced Global-Local Aggregation for Video Object Detection. CVPR 2020

Relation Networks for Object Detection (CVPR'2018)

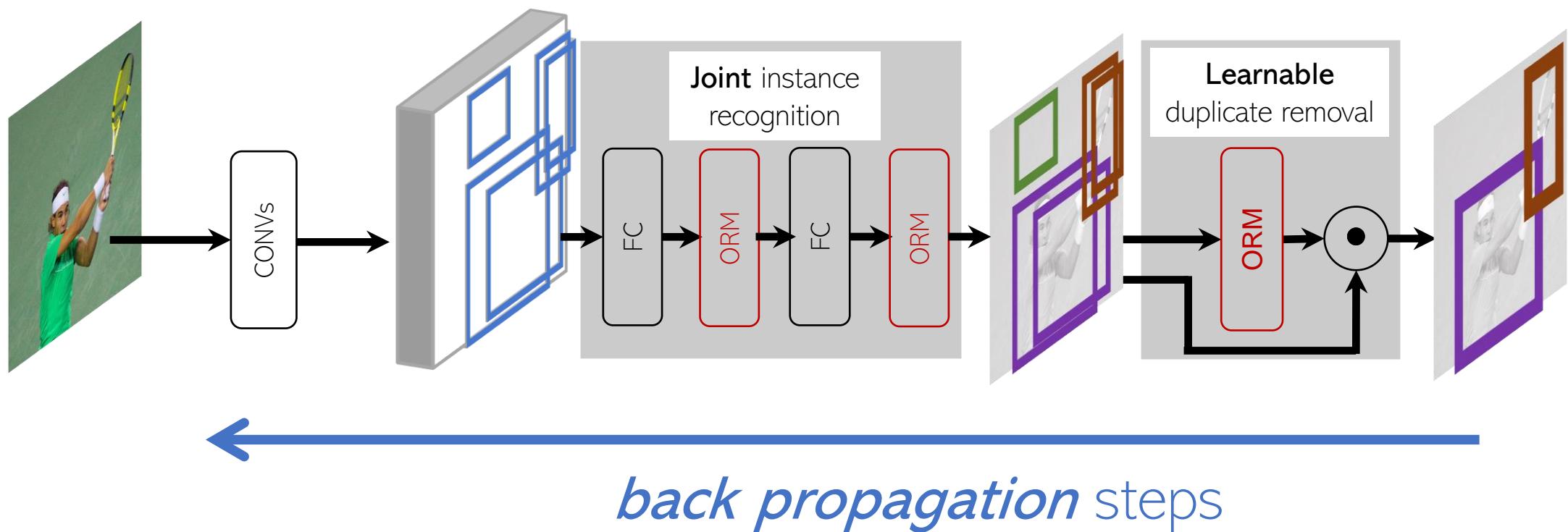


It is much easier to detect the *glove* if we know there is a *baseball player*.



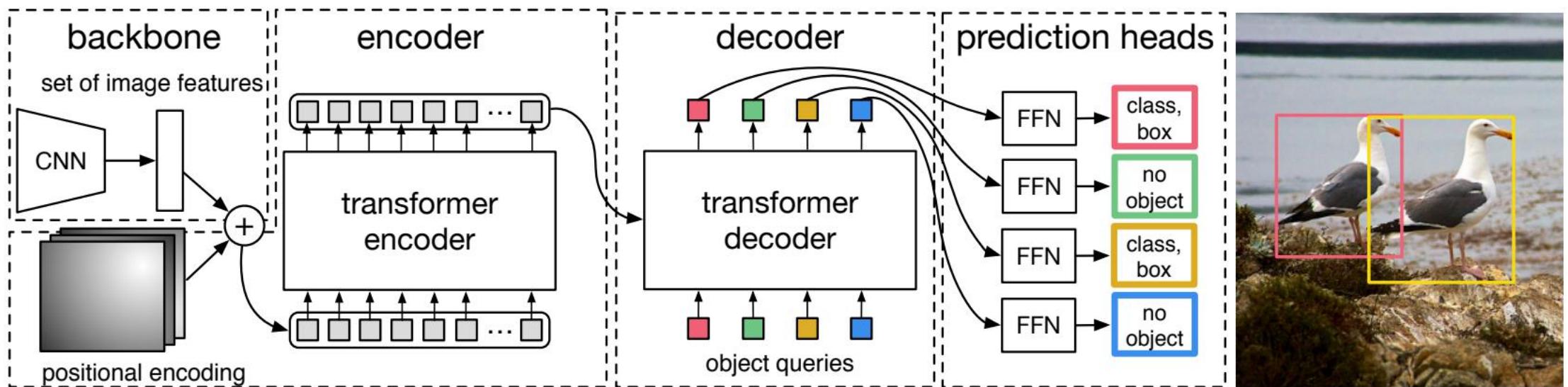
Relation Networks for Object Detection (CVPR'2018)

- The first fully end-to-end object detector



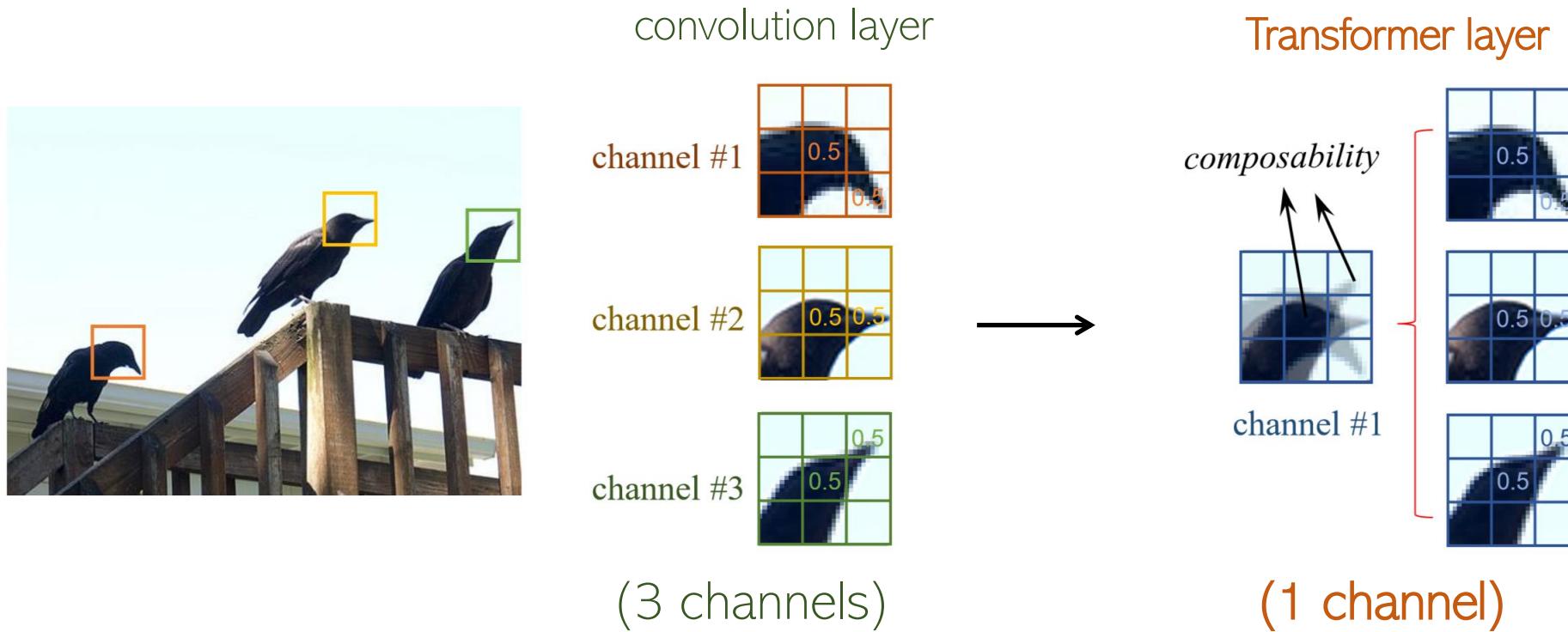
DETR (ECCV'2020)

- Another end-to-end object detector



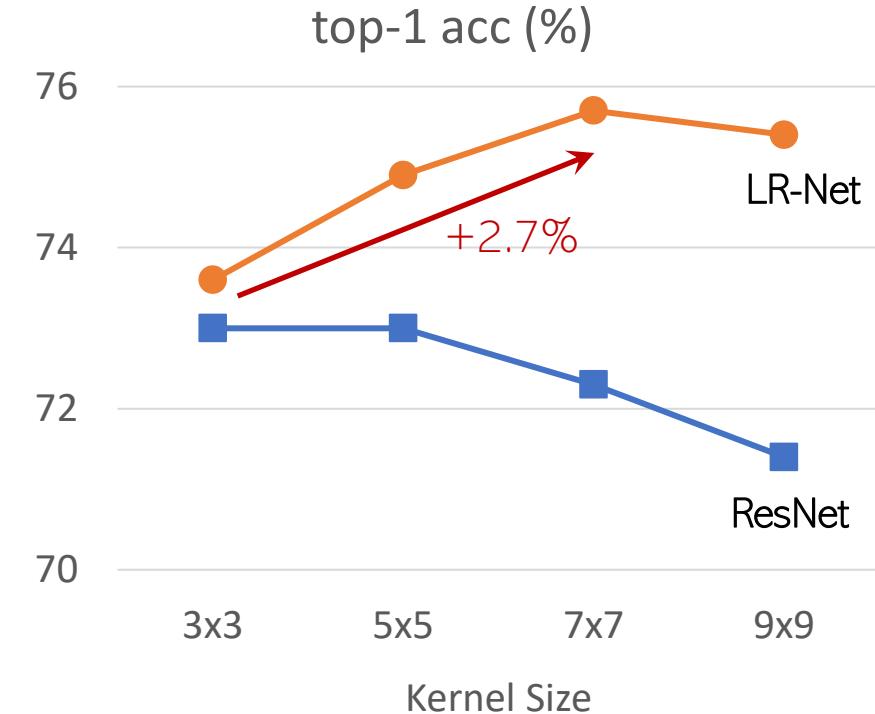
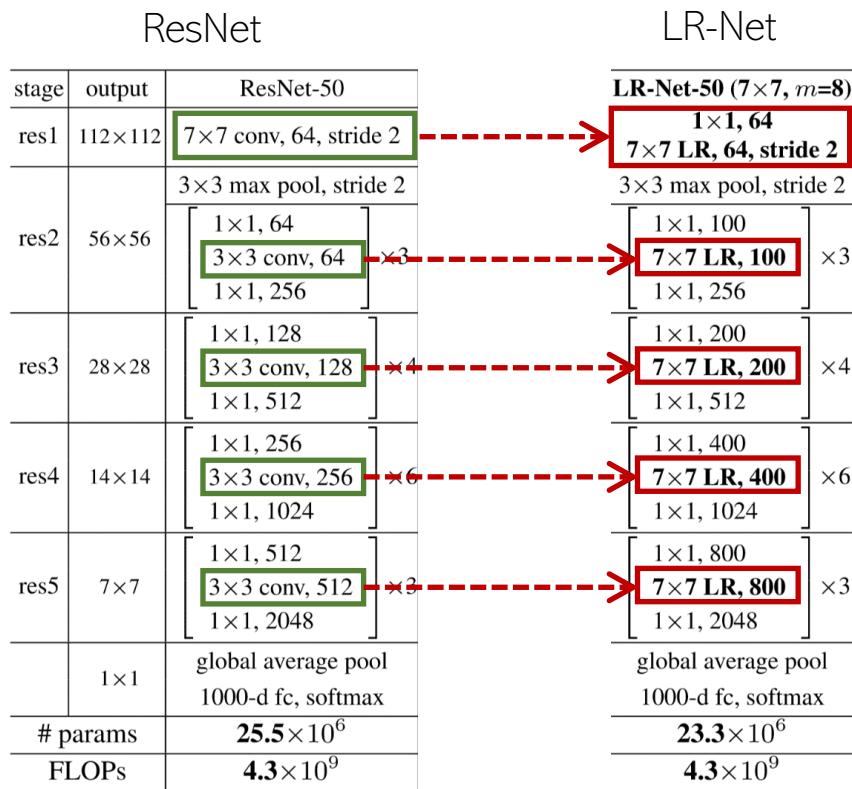
Local relation networks (ICCV'2019)

- Full-attention visual backbone
 - “Convolution is exponentially inefficient!”

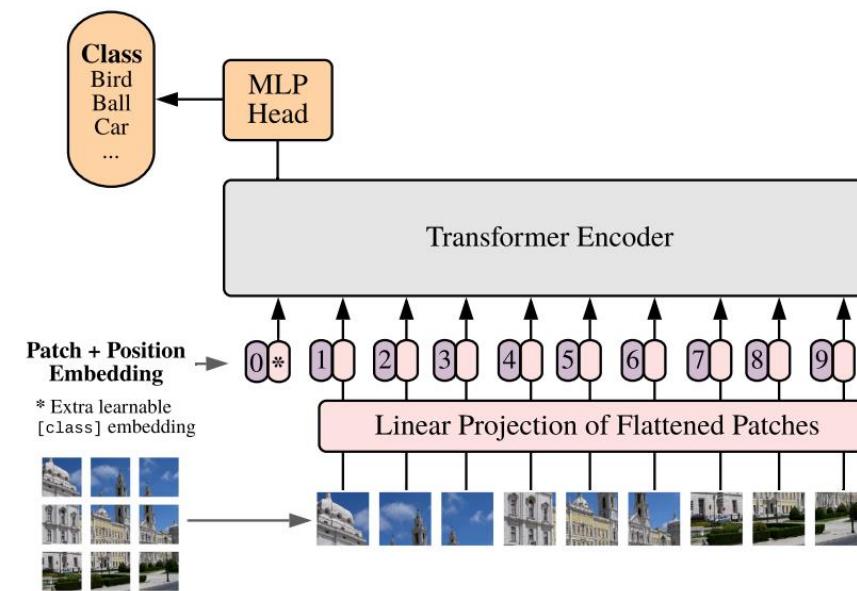
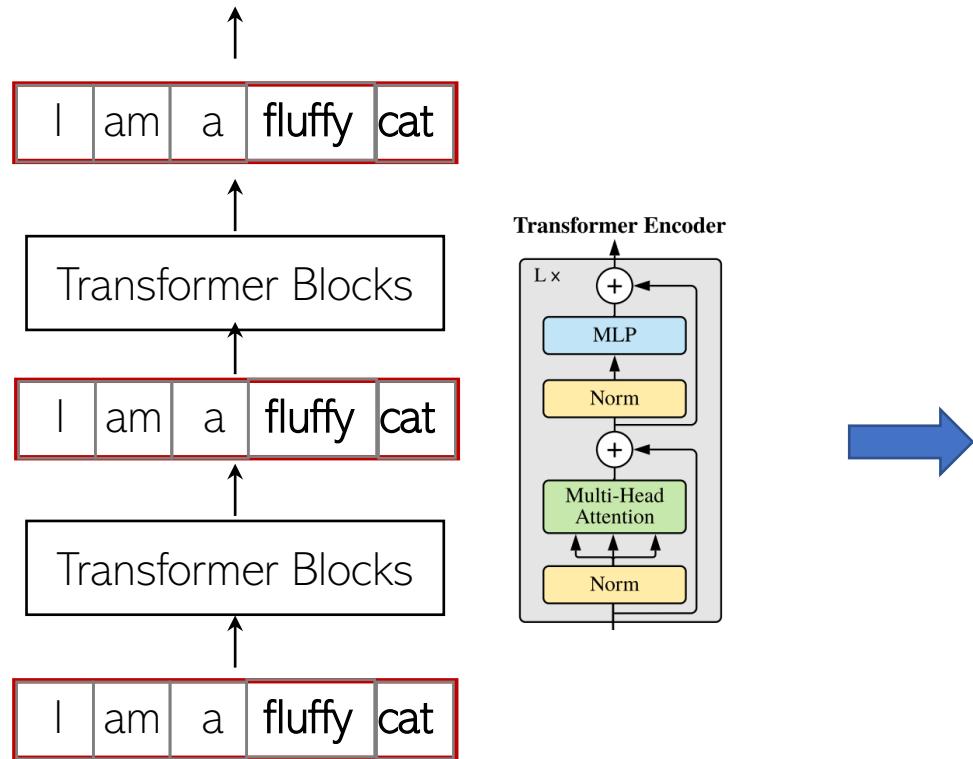


Local relation networks (ICCV'2019)

- The first full-attention visual backbone



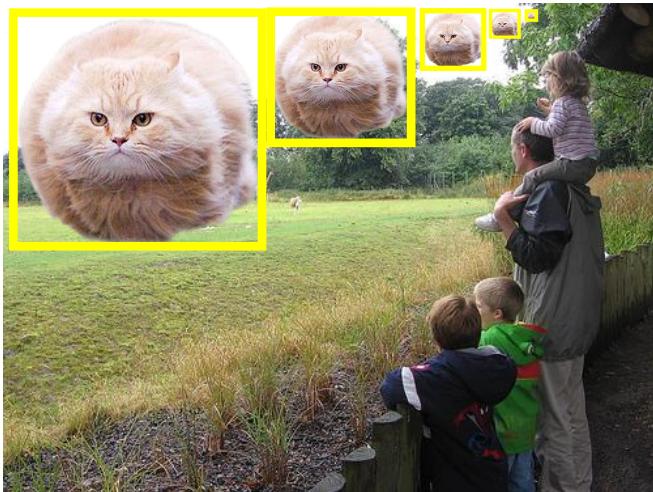
Vanilla Vision Transformer (ViT)



Key difference between visual and text signals

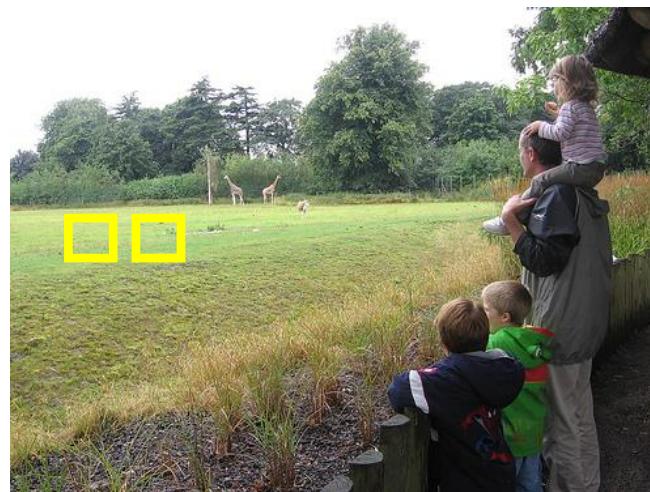
hierarchy

(scale invariance)

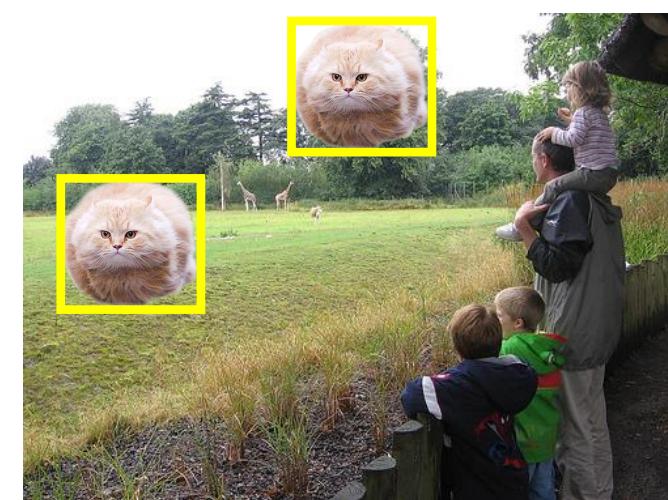


locality

(spatial smoothness)



translation invariance



I am a fluffy cat.

I am a fluffy fluffy cat cat. (invalid)

No scale variation

I like the green grass.

No spatial smoothness

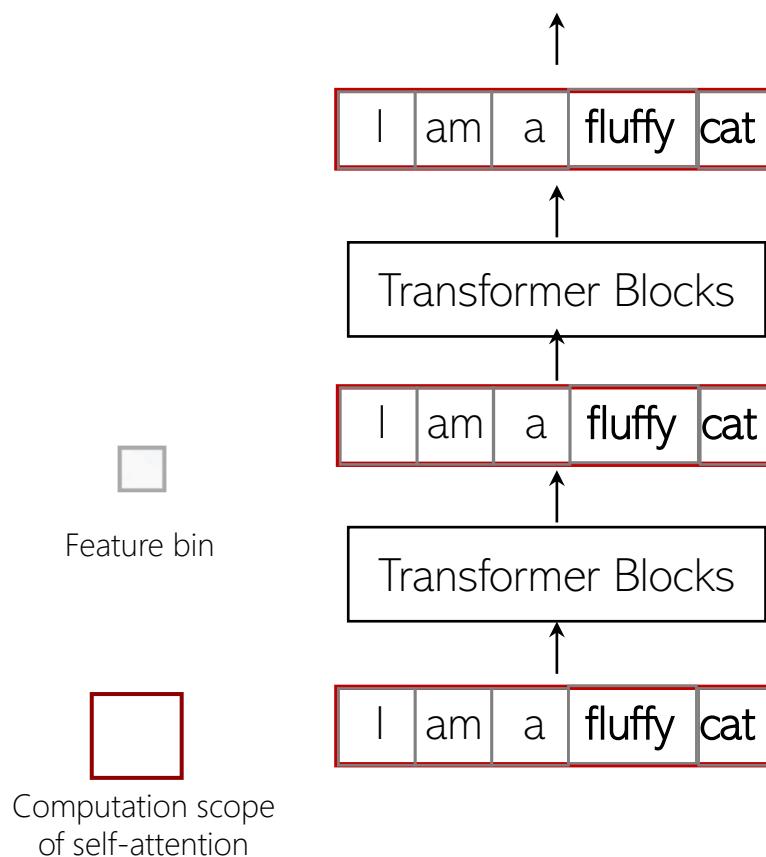
I am a fluffy cat.

Fluffy cat is me.

Sensitive to absolute locations

Computer vision priors: the bridge from NLP to CV

Standard Transformer (2017.6)

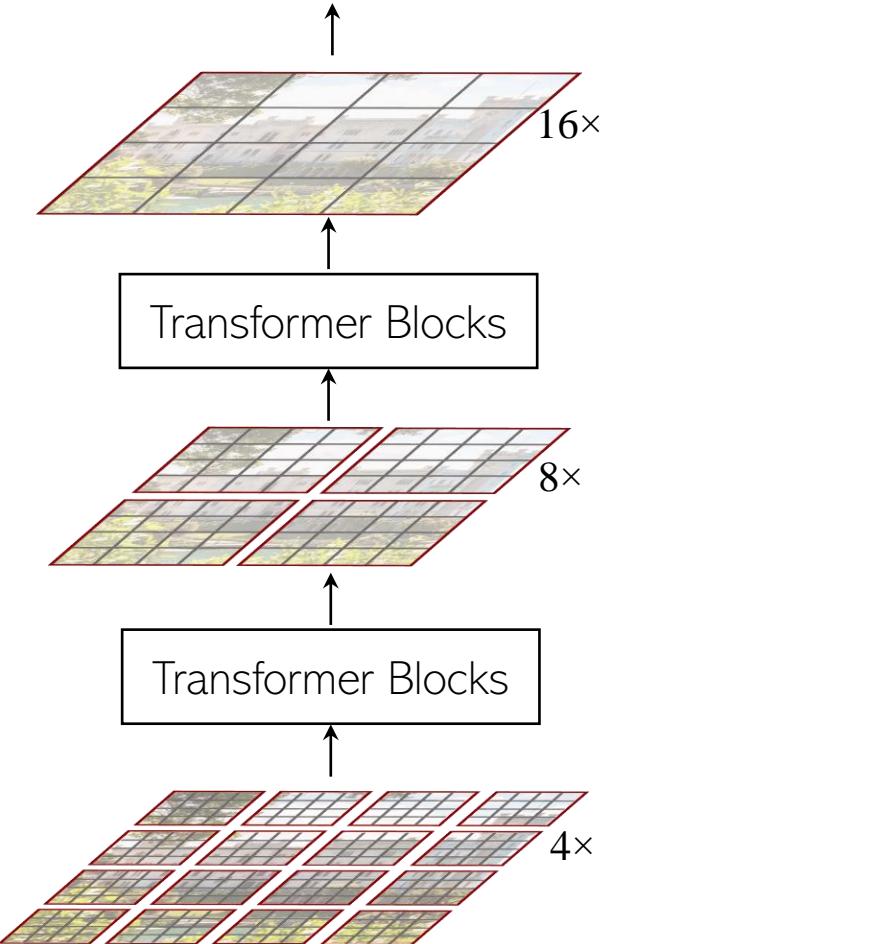


plain → hierarchy

global → locality

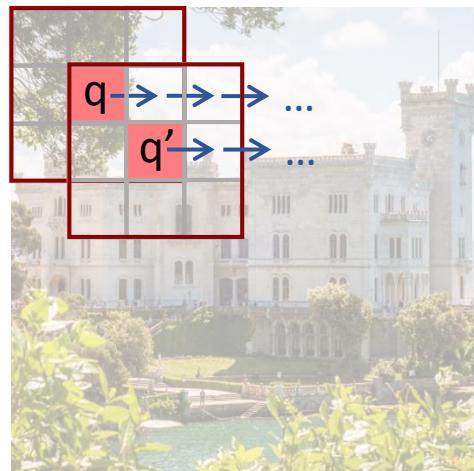
absolute position → translation invariance

Swin Transformer (2021.3)

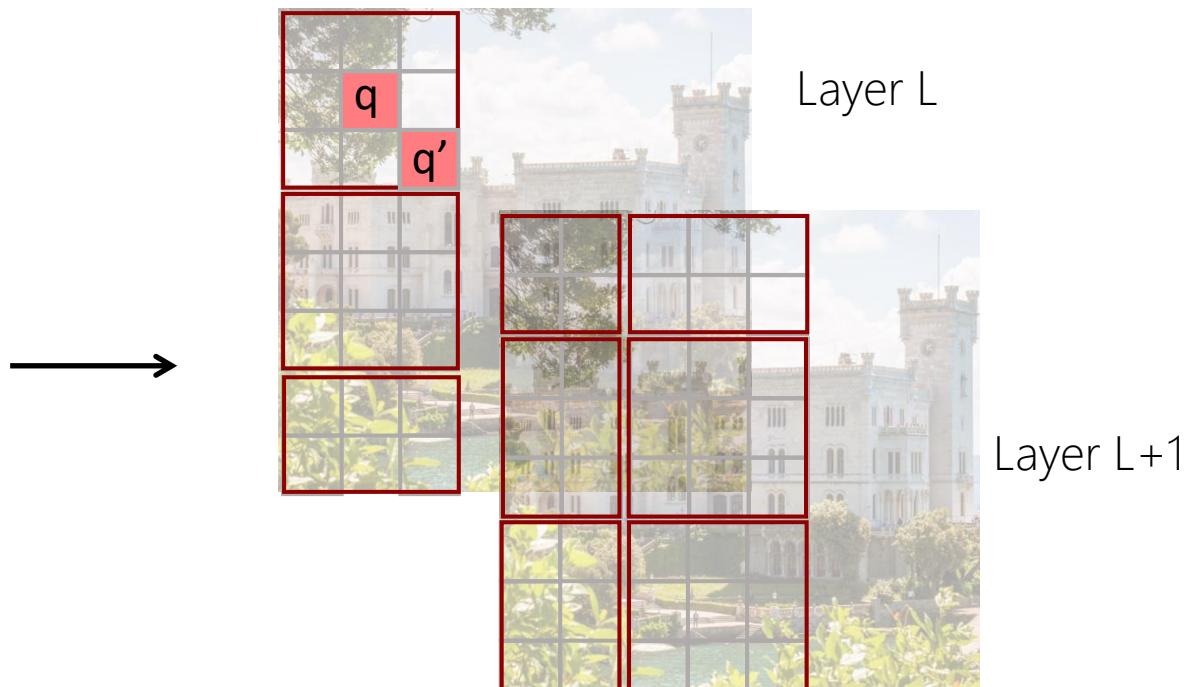


From sliding windows to shifted windows

- Non-overlapping windows (3x speed up in latency)
- Shifted window configurations in the next layer



Sliding window
(CNN/LR-Net)



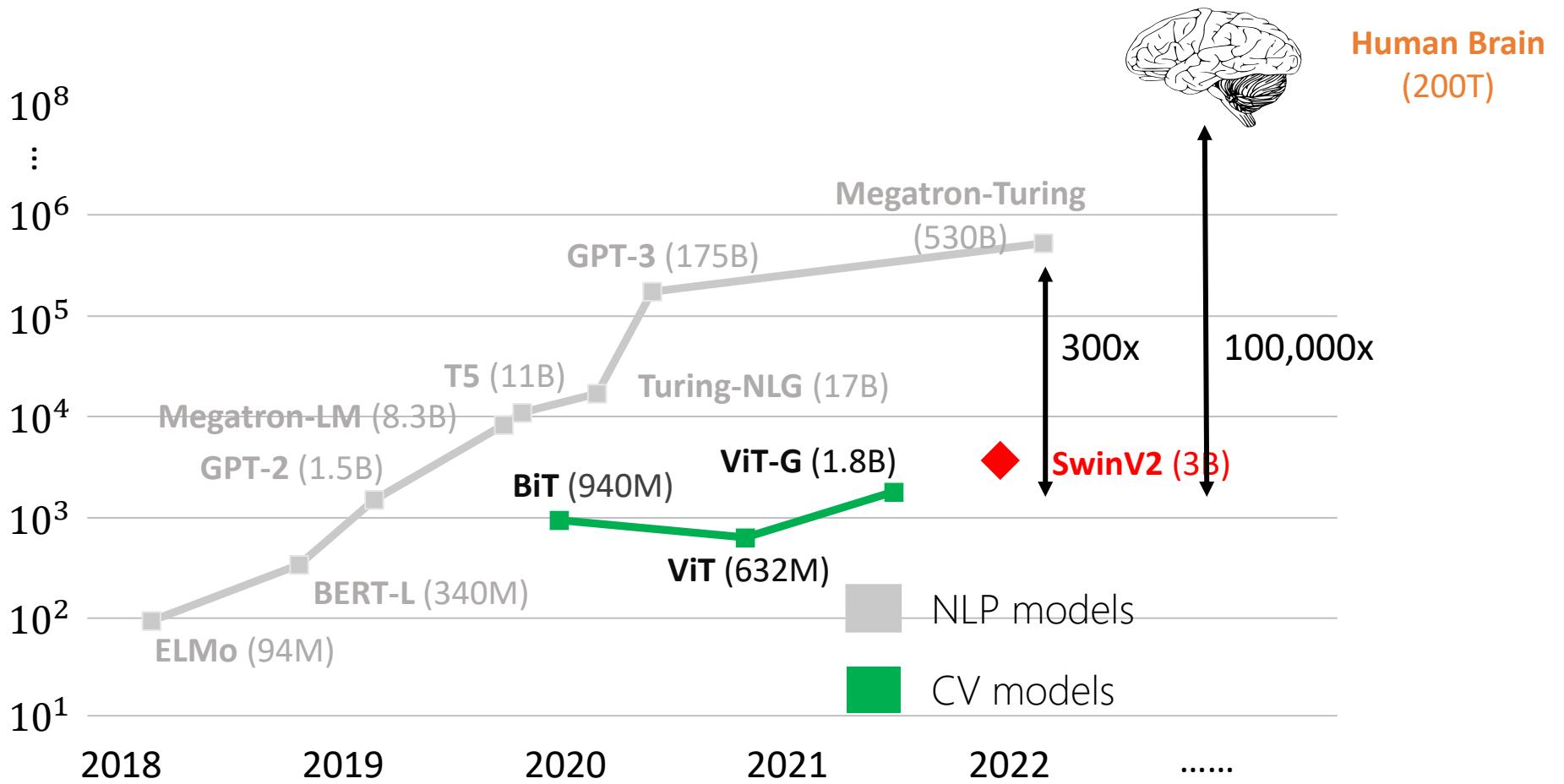
Shifted window (Swin, ICCV'2021)

Scaling up visual models

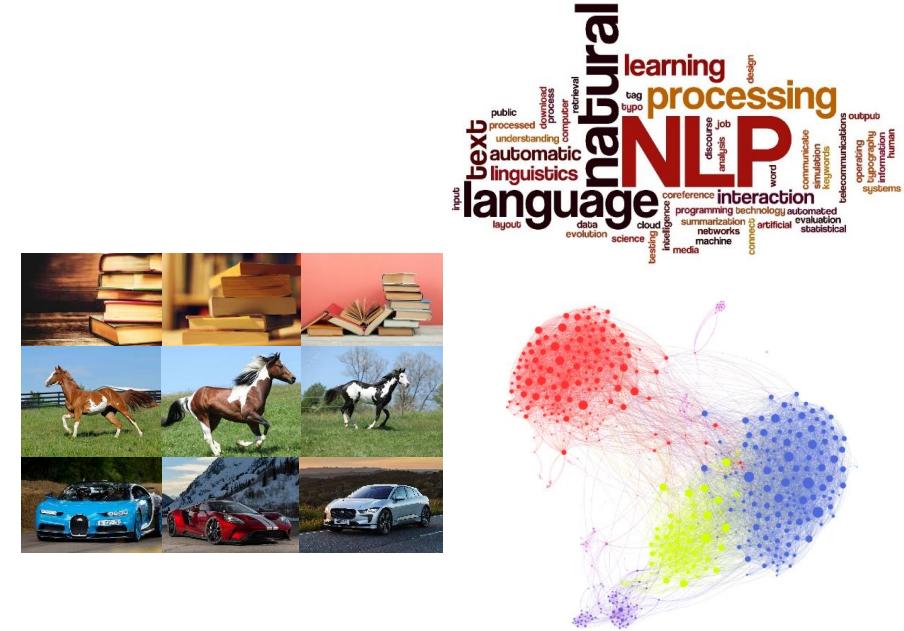
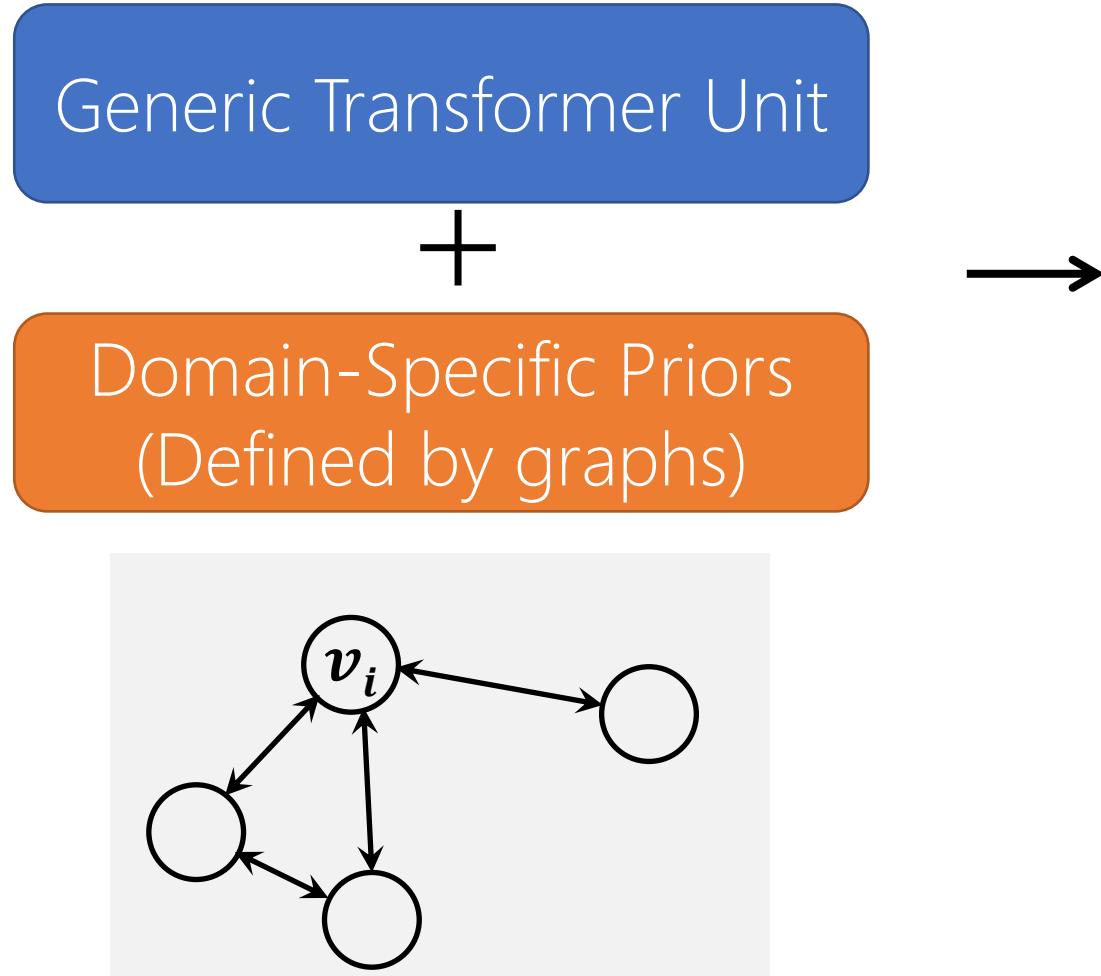
- Swin Transformer V2 (CVPR'2022)
 - The world's largest dense visual model (3B parameters)
 - Applicability of billion-scale visual models beyond image classification
 - Requiring 40 times less labelled data than Google's

Benchmark	ImageNetV2	COCO test-dev	ADE20K val	Kinetics-400
Swin V1	77.5	58.7	53.5	84.9
Previous SoTA	83.3 (Google, July 2021)	61.3 (Microsoft, July 2021)	58.4 (Microsoft, October 2021)	85.4 (Google, October 2021)
Swin V2 (Giant) (November 2021)	84.0 (+0.7)	63.1 (+1.8)	59.9 (+1.5)	86.8 (+1.4)

Computer vision models are still relatively small



A general principle to apply Transformer



OGB-LSC
— Large - Scale Challenge —

Dataset: Predicting a quantum property of molecular graphs

Graphomer, NeurIPS 2021

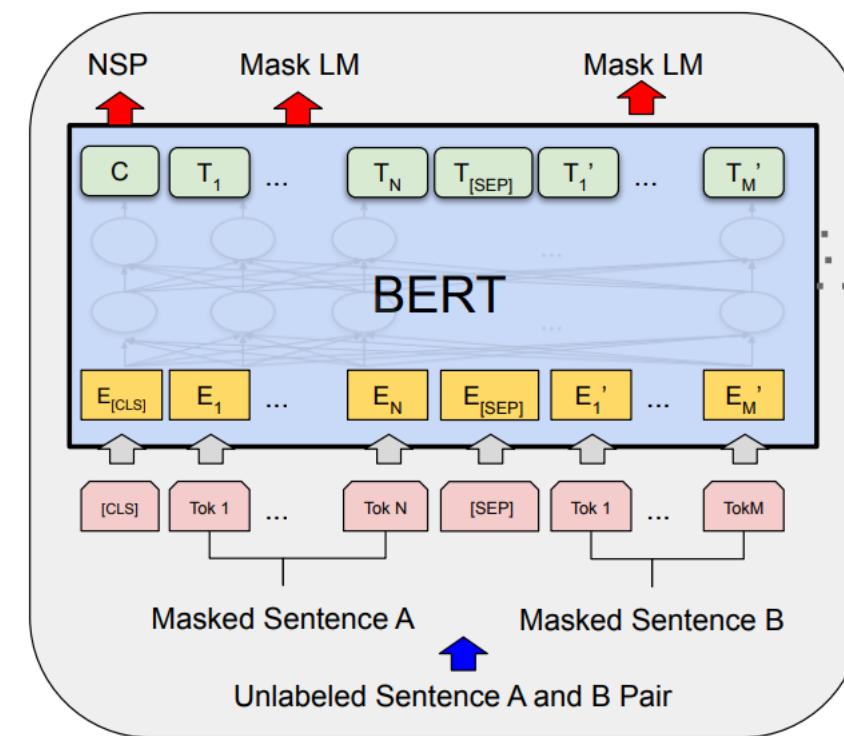
Towards converged learning methods

Masked signal modeling (BERT-like)

Mainstream pre-training methods for AI subfields

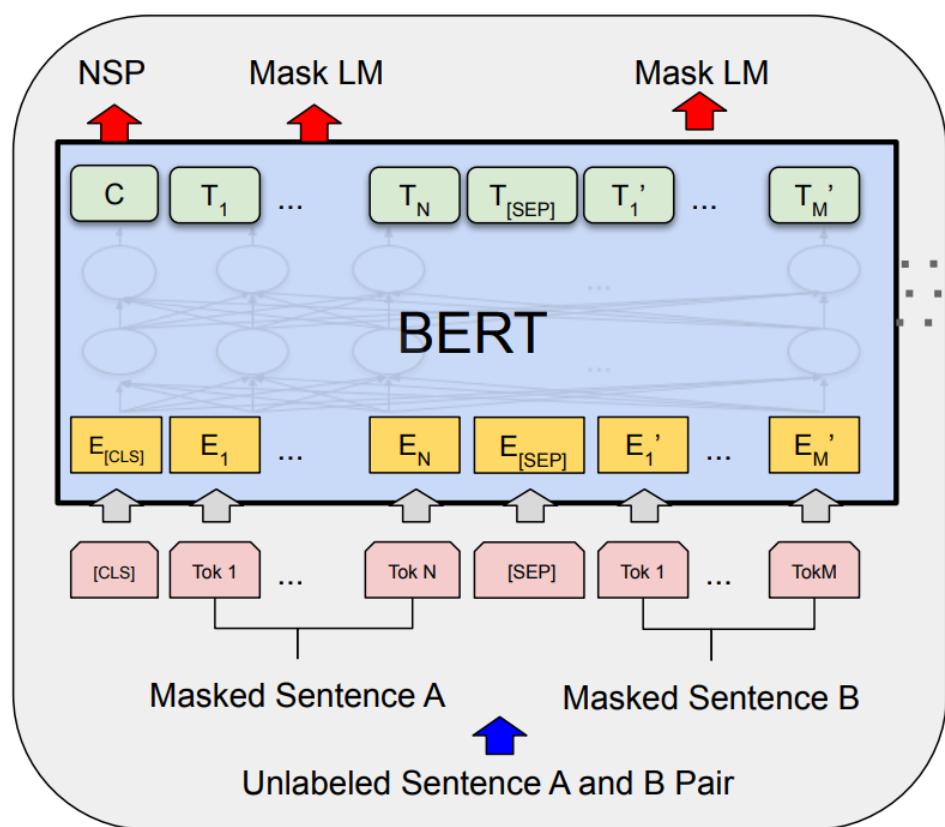


Supervised learning
(Computer vision)

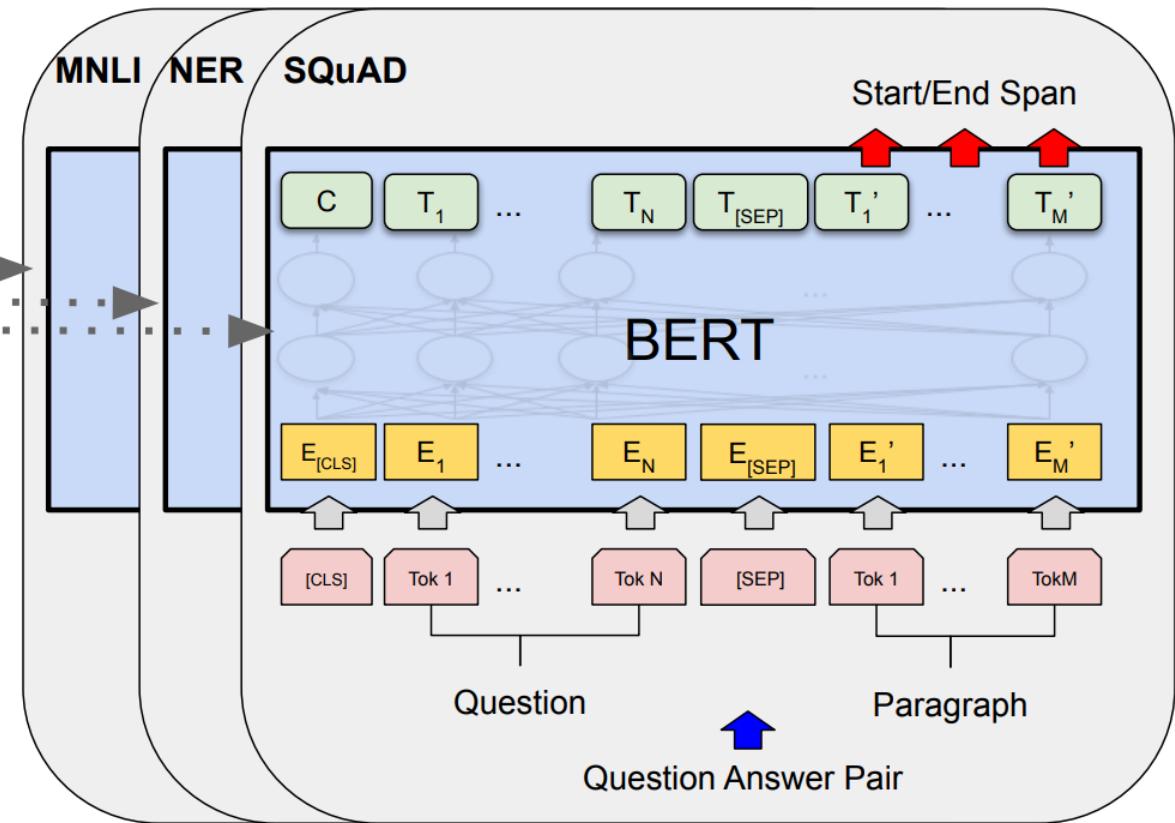


Self-supervised learning
(Nature language processing)

BERT in NLP: Predicting unknown input



Pre-training



Fine-Tuning

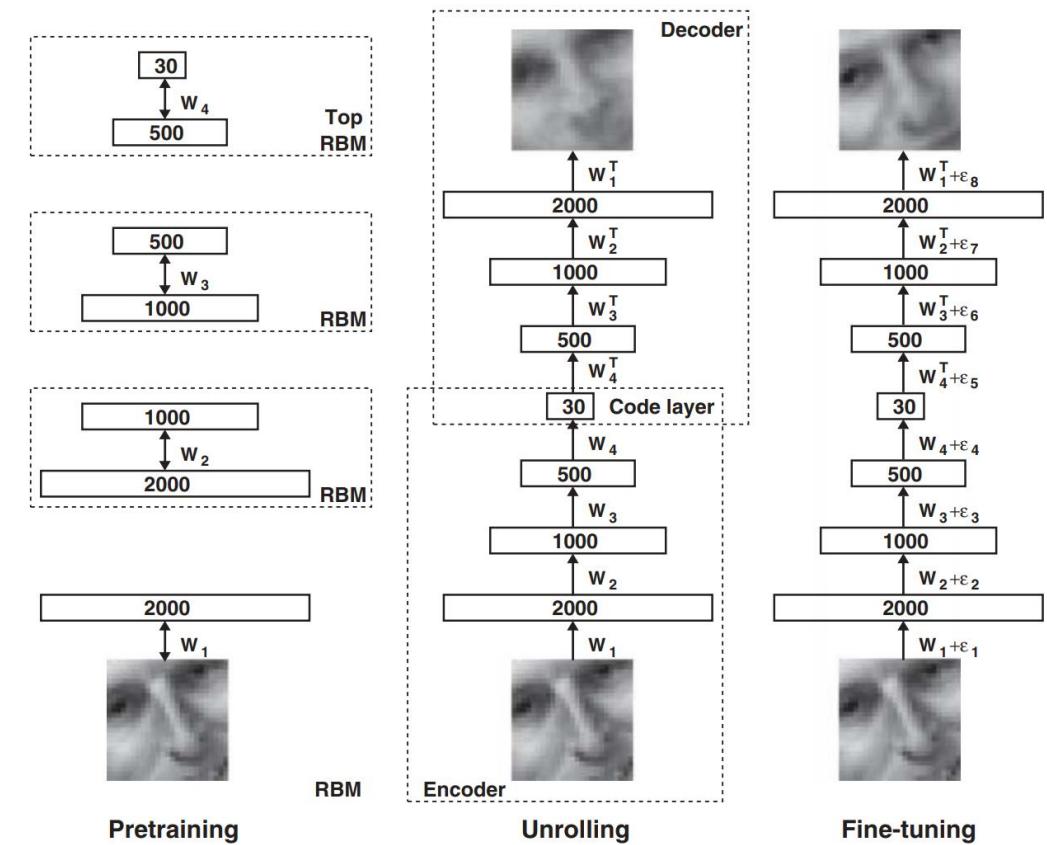
Self-supervised learning in computer vision

- Pioneer of deep learning

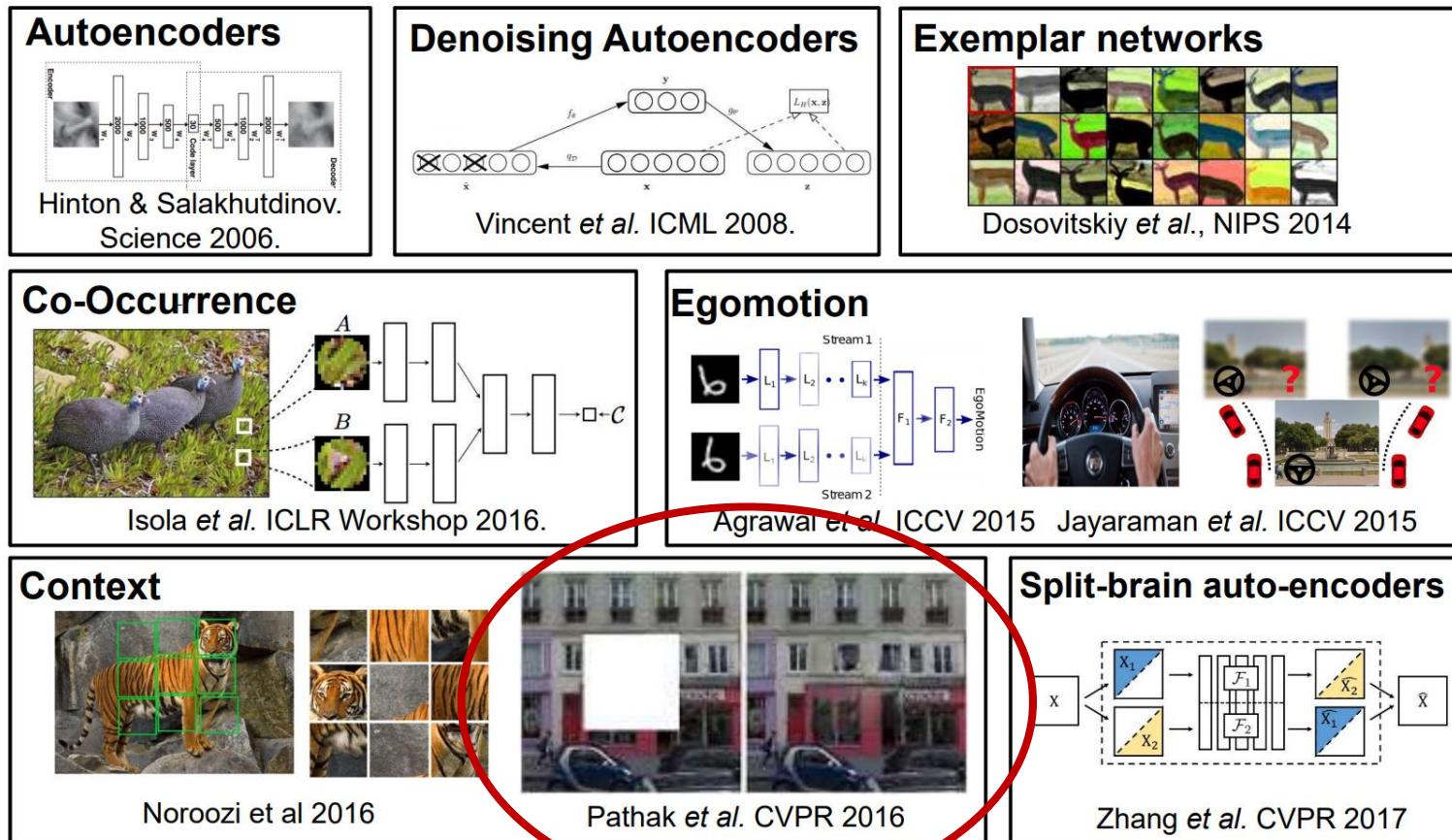
Science, 2006

Reducing the Dimensionality of Data with Neural Networks

G. E. Hinton* and R. R. Salakhutdinov



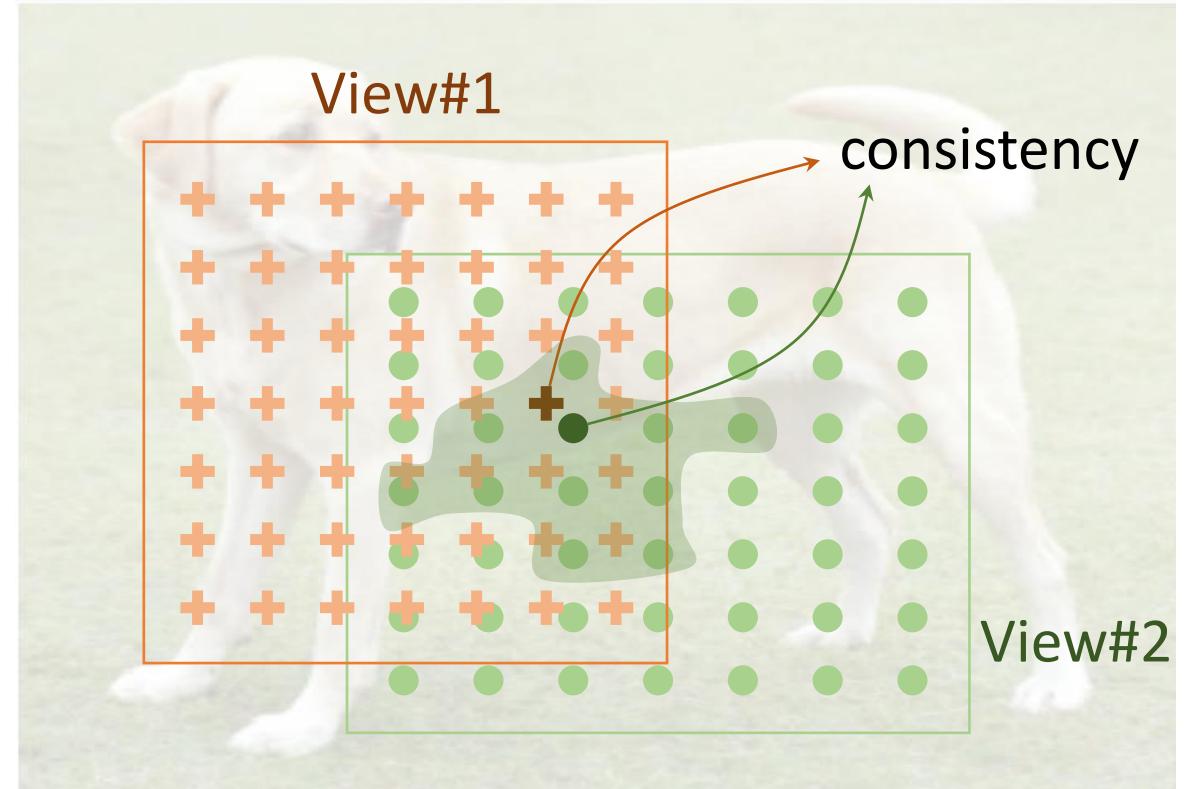
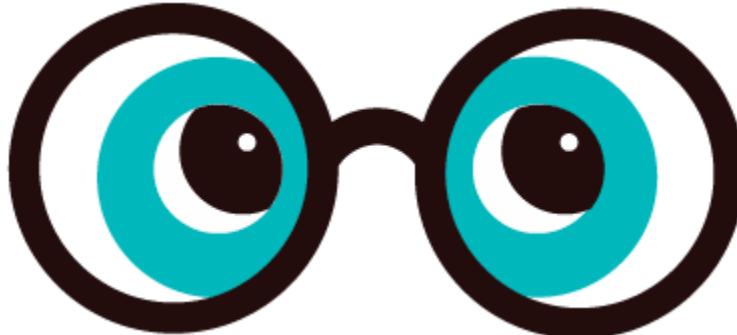
Self-supervised learning in computer vision



From Andrew Zisserman

Self-supervised learning in CV: Prediction of unknown input or using contrast

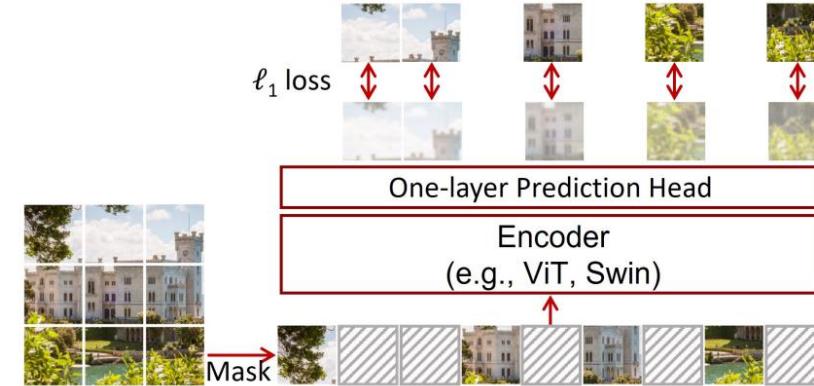
- MoCo (CVPR'2020)
- PixPro (CVPR'2021)
 - Inspired by eye movements



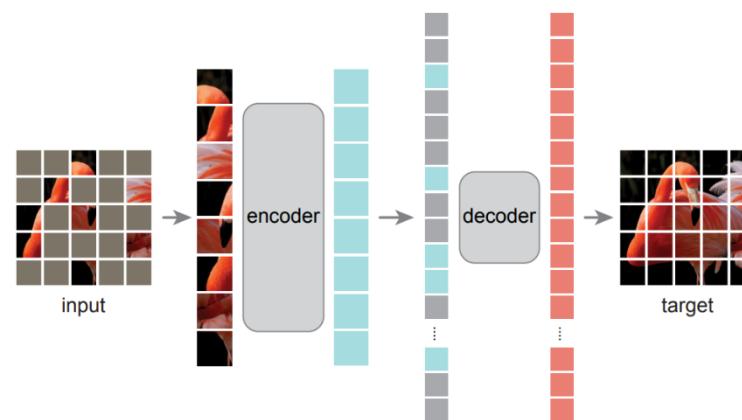
PixPro method

Self-supervised learning in CV: Prediction of unknown input

- GPT/BERT-like approaches
 - Image-GPT (OpenAI)
 - BEiT (MSRA)
 - **SimMIM (MSRA)**
 - MAE (Meta)
 - PeCo (MSRA)
 - data2vec (Meta)
 - CAE (Baidu)
 - ...



SimMIM: A Simple Framework for Masked Image Modeling



MAE: Masked Autoencoders Are Scalable Vision Learners

BERT-like approach: differences between CV and NLP

- SimMIM (CVPR'2022)

I am a **fluffy** cat

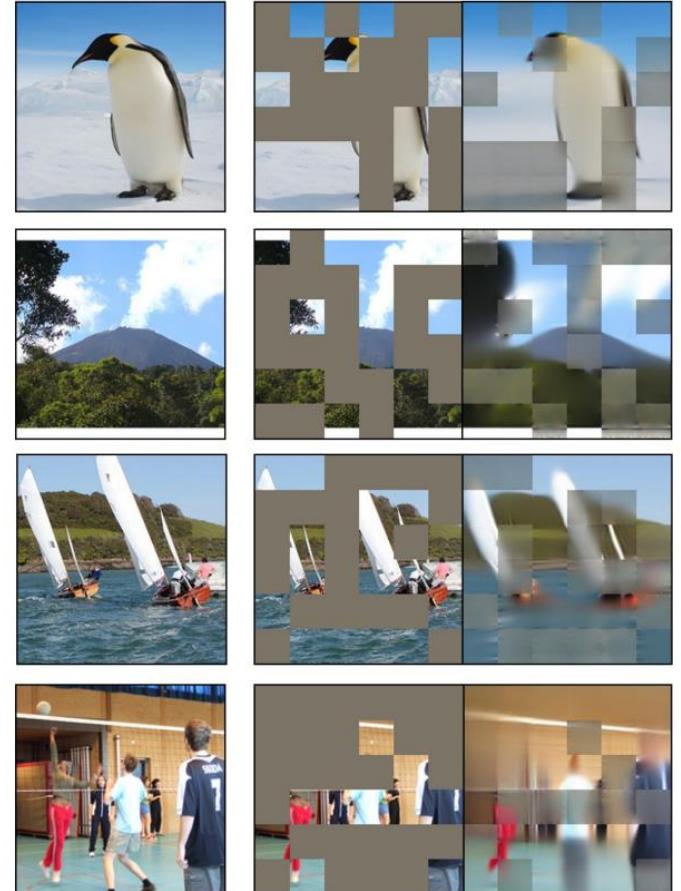


I am a [mask] cat

- ✓ No spatial smoothness
- ✓ High-level discrete



- ✓ Spatial smoothness
- ✓ Raw and low-level continuous

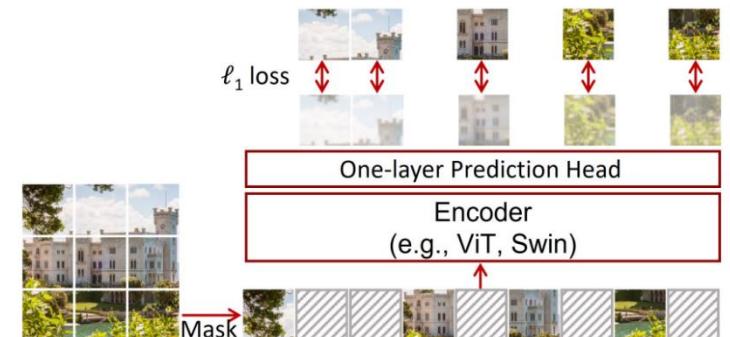
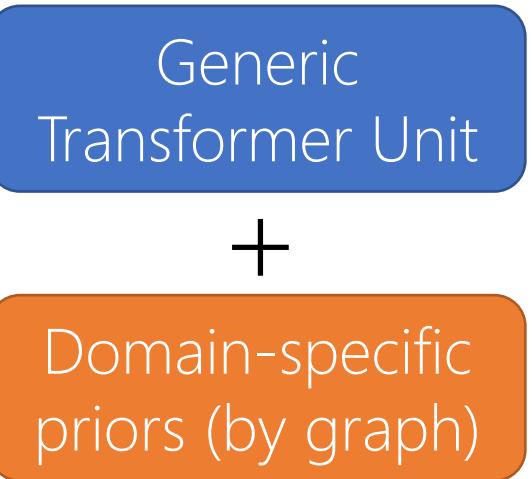


Low mask ratio (15%), token classification

Large patch size and mask ratio (ps=32, ratio=60%), direct pixel value regression

Take-home message

- Embrace converged architectures and learning methods in AI
 - Converged architectures
 - Transformer + domain-specific priors
 - NLP: Transformer
 - CV: **Swin Transformer V1/V2**
 - Converged pre-training methods
 - Prediction of unknown input
 - NLP: BERT/GPT/XLNet
 - CV: ImageGPT/BEiT/**SimMIM**/MAE/PeCo/data2vec/CAE ...



Acknowledgements

- Swin Transformer series benefit a lot from OpenMMLab

This is an official implementation for "Swin Transformer: Hierarchical Vision Transformer using Shifted Windows".

91c54bf · 12 days ago · 59 commits

- configs add Swin-MoE (#220) 15 days ago
- data add Swin-MoE (#220) 15 days ago
- figures Initial commit 14 months ago
- models add Swin-MoE (#220) 15 days ago
- .gitignore Initial commit 14 months ago
- CODE_OF_CONDUCT.md Initial CODE_OF_CONDUCT.md commit 15 months ago
- LICENSE Updating LICENSE to template content 15 months ago
- MODELHUB.md add Swin-MoE (#220) 15 days ago
- README.md Update README.md 12 days ago
- SECURITY.md Initial SECURITY.md commit 15 months ago
- SUPPORT.md Initial SUPPORT.md commit 15 months ago
- config.py add Swin-MoE (#220) 15 days ago
- get_started.md add Swin-MoE (#220) 15 days ago
- logger.py Initial commit 14 months ago
- lr_scheduler.py Initial commit 14 months ago
- main.py Change Apex amp to Pytorch amp (#207) last month
- main_moe.py add Swin-MoE (#220) 15 days ago
- optimizer.py Initial commit 14 months ago
- utils.py add imagenet22k dataset and some minor fixes (#208) last month
- utils_moe.py add Swin-MoE (#220) 15 days ago

README.md

Swin Transformer

Object Detection on COCO test dev (using additional training data)
Semantic Segmentation on COCO2014 val (using additional training data)
Action Recognition on K400 (using additional training data)
Action Classification on Kinetics-400 (using additional training data)

By Ze Liu*, Yutong Lin*, Yue Cao*, Han Hu*, Yixuan Wei, Zheng Zhang, Stephen Lin and Baining Guo.

This repo is the official implementation of "Swin Transformer: Hierarchical Vision Transformer using Shifted Windows". It currently includes code and models for the following tasks:

- Image Classification: Included in this repo. See `get_started.md` for a quick start.
- Object Detection and Instance Segmentation: See `Swin Transformer for Object Detection`.
- Semantic Segmentation: See `Swin Transformer for Semantic Segmentation`.
- Video Action Recognition: See `Video Swin Transformer`.
- Semi-Supervised Object Detection: See `Soft Teacher`.
- SSL: Contrastive Learning: See `Transformer-SSL`.
- SSL: Masked Image Modeling: See `SimMIM`.



Swin Transformer

This organization maintains repositories built on Swin Transformers. The pretrained models locate at <https://github.com/microsoft/Swin-Transformer>

<https://arxiv.org/pdf/2103.14030.pdf>

Follow

Overview Repositories 6 Projects Packages Teams People 7 Settings

Pinned

Customize pins

Transformer-SSL Public Forked from microsoft/Swin-Transformer This is an official implementation for "Self-Supervised Learning with Swin Transformers". Python ⭐ 487 📂 55

Swin-Transformer-Object-Detection Public Forked from open-mmlab/mmldetection This is an official implementation for "Swin Transformer: Hierarchical Vision Transformer using Shifted Windows" on Object Detection and Instance Segmentation. Python ⭐ 1.2k 📂 302

Swin-Transformer-Semantic-Segmentation Public Forked from open-mmlab/mmsegmentation This is an official implementation for "Swin Transformer: Hierarchical Vision Transformer using Shifted Windows" on Semantic Segmentation. Python ⭐ 788 📂 178

Video-Swin-Transformer Public Forked from open-mmlab/mmaction2 This is an official implementation for "Video Swin Transformers". Python ⭐ 787 📂 125

View as: Public You are viewing this page as a public user. You can [create a README file](#) visible to anyone.

People

Invite someone

Top languages

Python