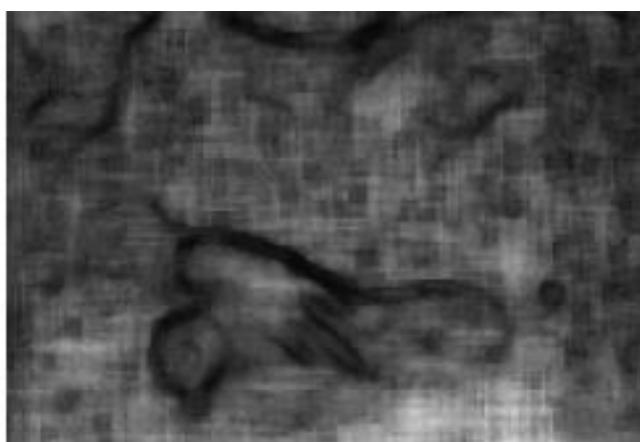
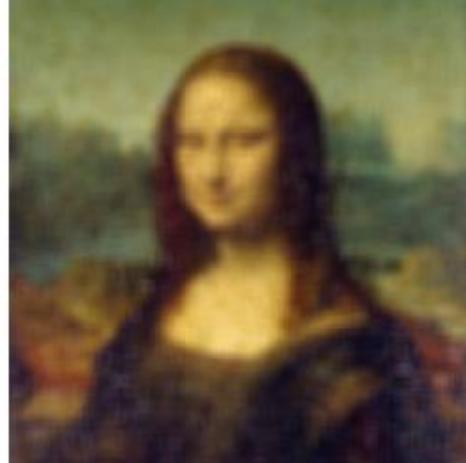


# Computer Vision and Image Processing

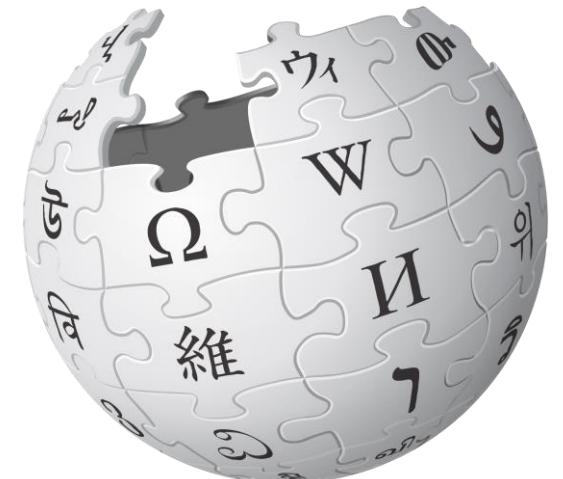
[Sanparith.Marukat@nectec.or.th](mailto:Sanparith.Marukat@nectec.or.th)

IPU / AINRG / NECTEC | AIAT



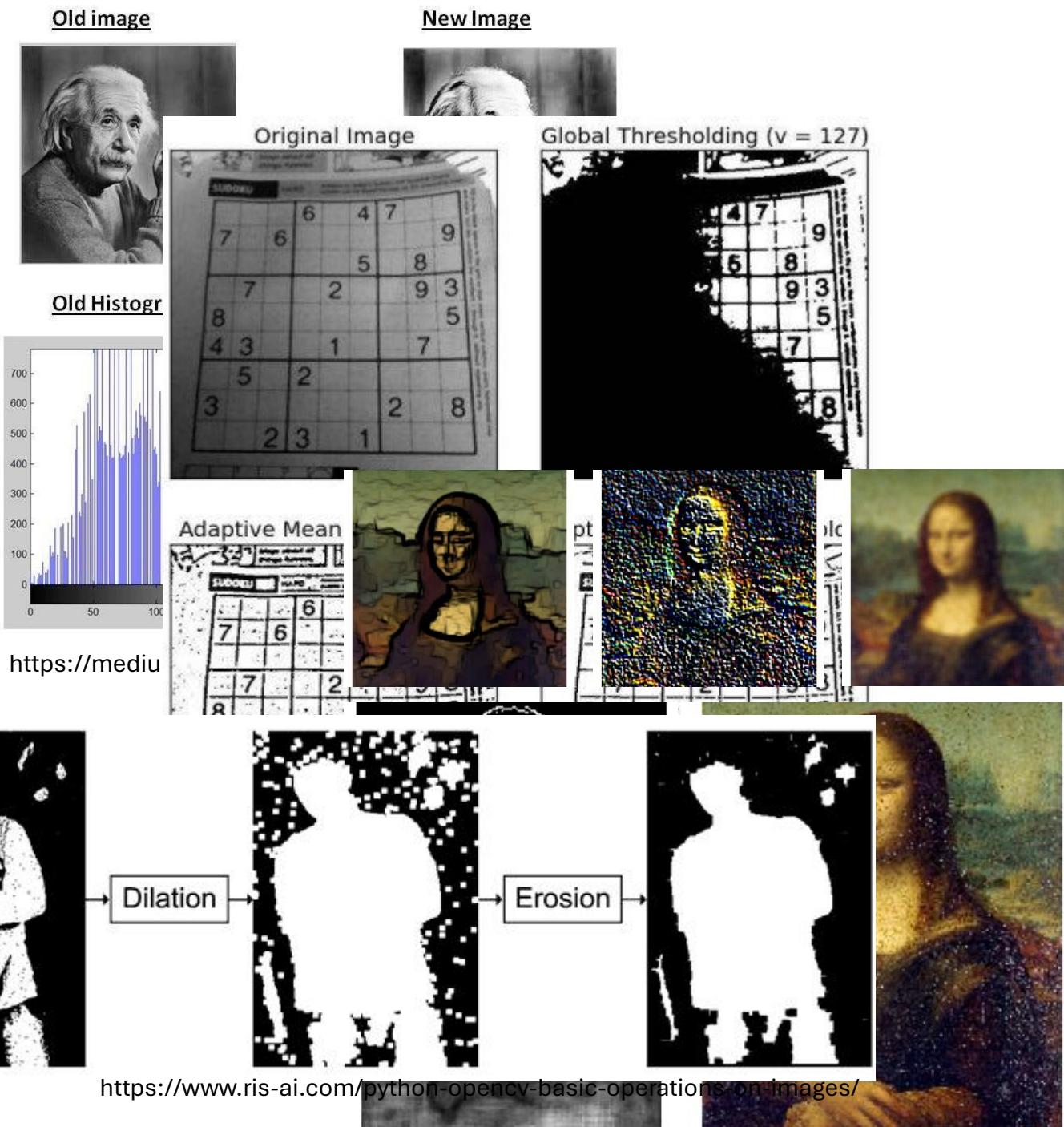
# Definition

- An interdisciplinary field that deals with how computers can be made to gain **high-level understanding** from **digital images or videos**.
- From the perspective of engineering, it seeks to automate **tasks** that the human visual system can do.



# Tasks (low level)

- Intensity transformation
  - Histogram equalization
  - Contrast stretching
- Image binarization (thresholding)
- Image filtering
  - Smoothing
  - Edge detection
  - Sharpening
- Morphological operations
  - Erosion and Dilation
  - Opening and Closing

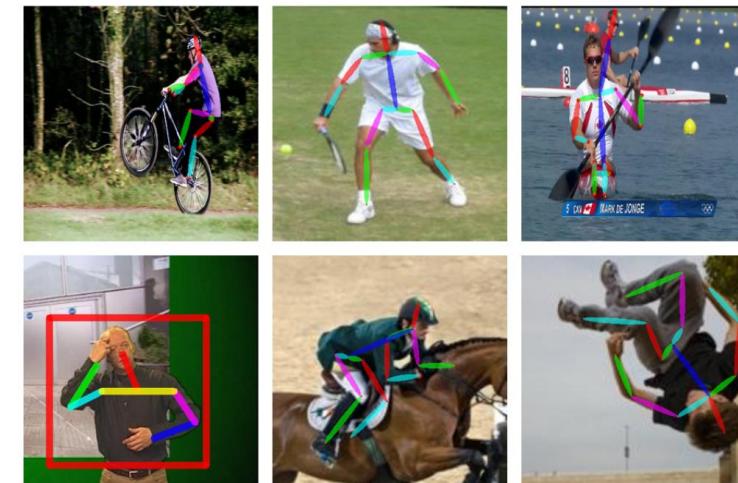
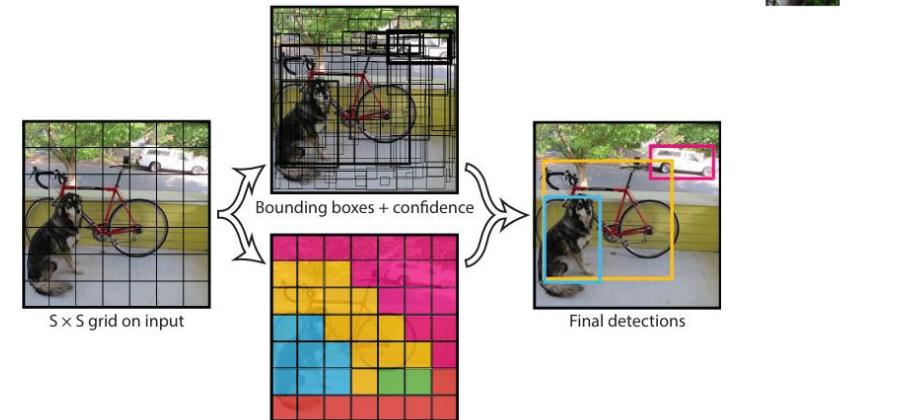
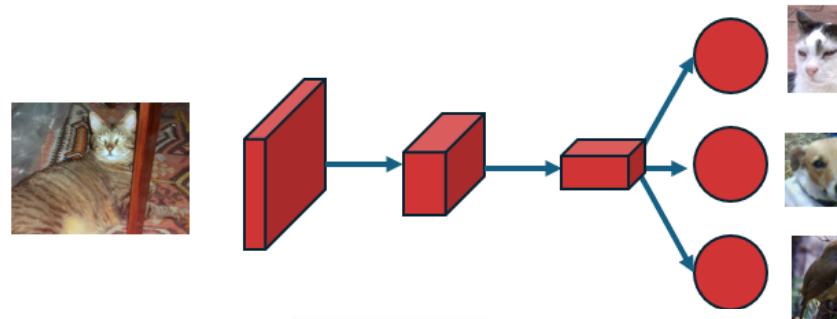


# Tasks (high level)

- Image classification
- Object detection
  - Tracking
  - Pose estimation



<https://opencv.org/blog/multiple-object-tracking-in-realtime/>



<https://arxiv.labs.arxiv.org/html/2001.08>

# Tasks (high level)

- Image segmentation
  - Semantic
  - Instance
  - Panoptic
- Depth estimation
- Structure from motion

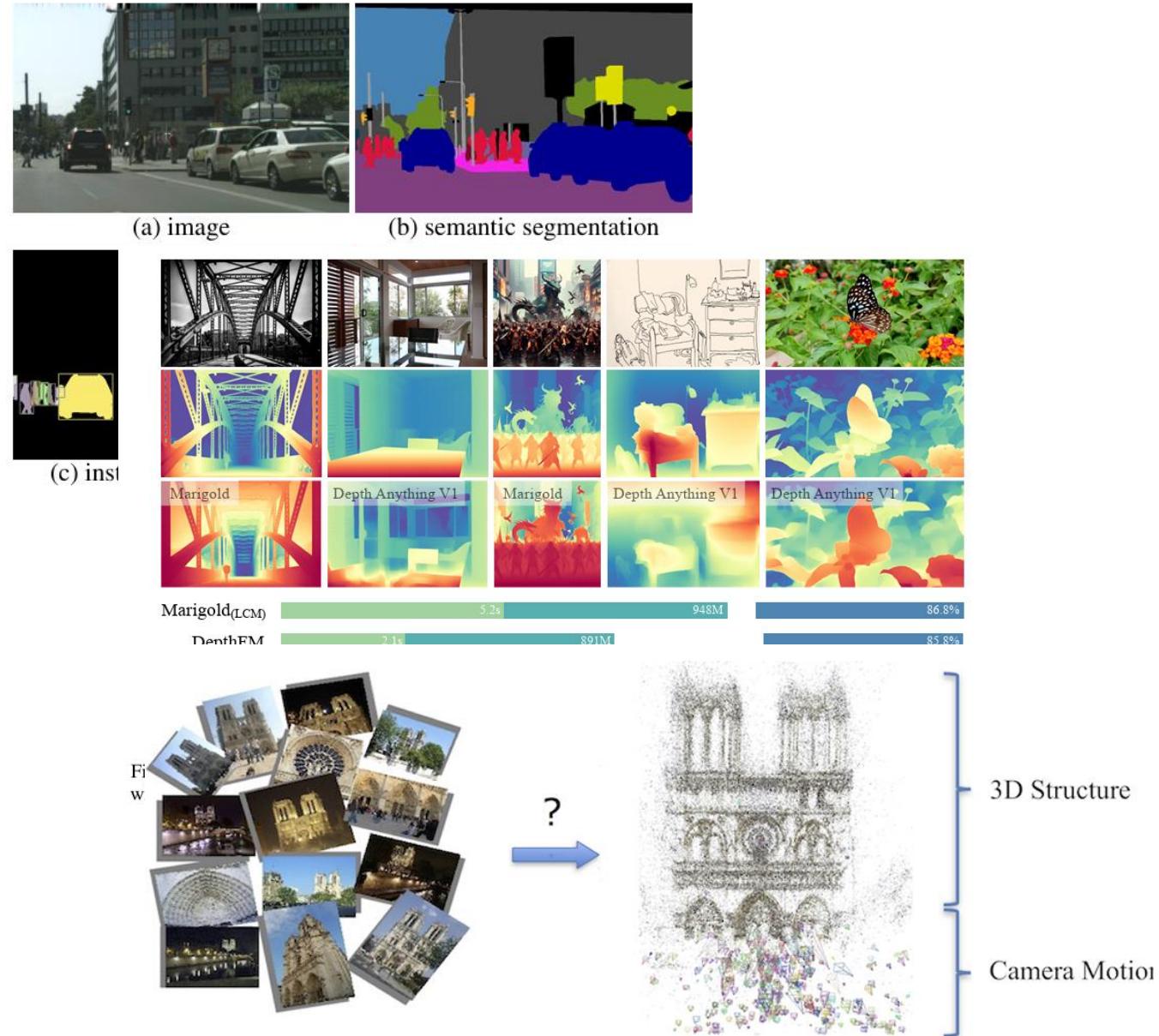


Figure 1: The structure from motion (SfM) problem.

<https://arxiv.org/abs/1701.08493>

# Tasks (high level)

- Image Search
  - Content-Based Image Retrieval
  - Object Search
  - Image Search from text query
- Multi-Modal
  - Image Captioning
  - Visual Question Answering
  - Image Generation from text prompt

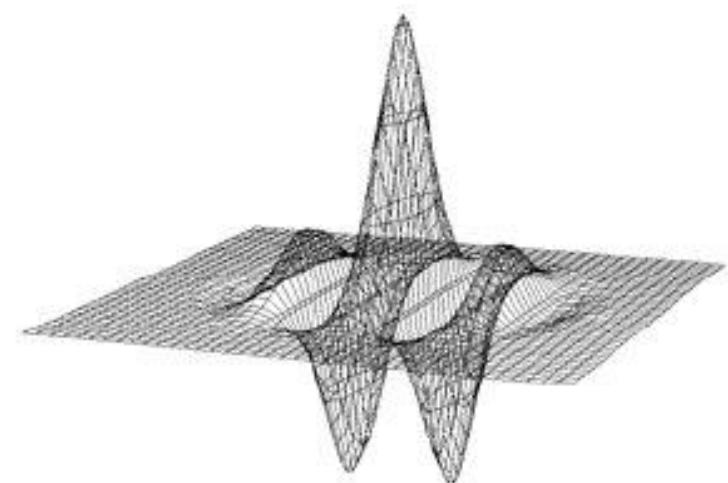
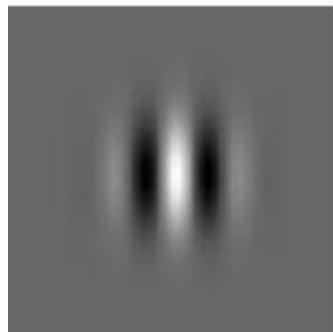
# Tasks (More advanced)

- Training framework
  - Self-supervised learning
  - Imbalanced dataset
  - Labelling error
- Deployment: GPU?

# Techniques

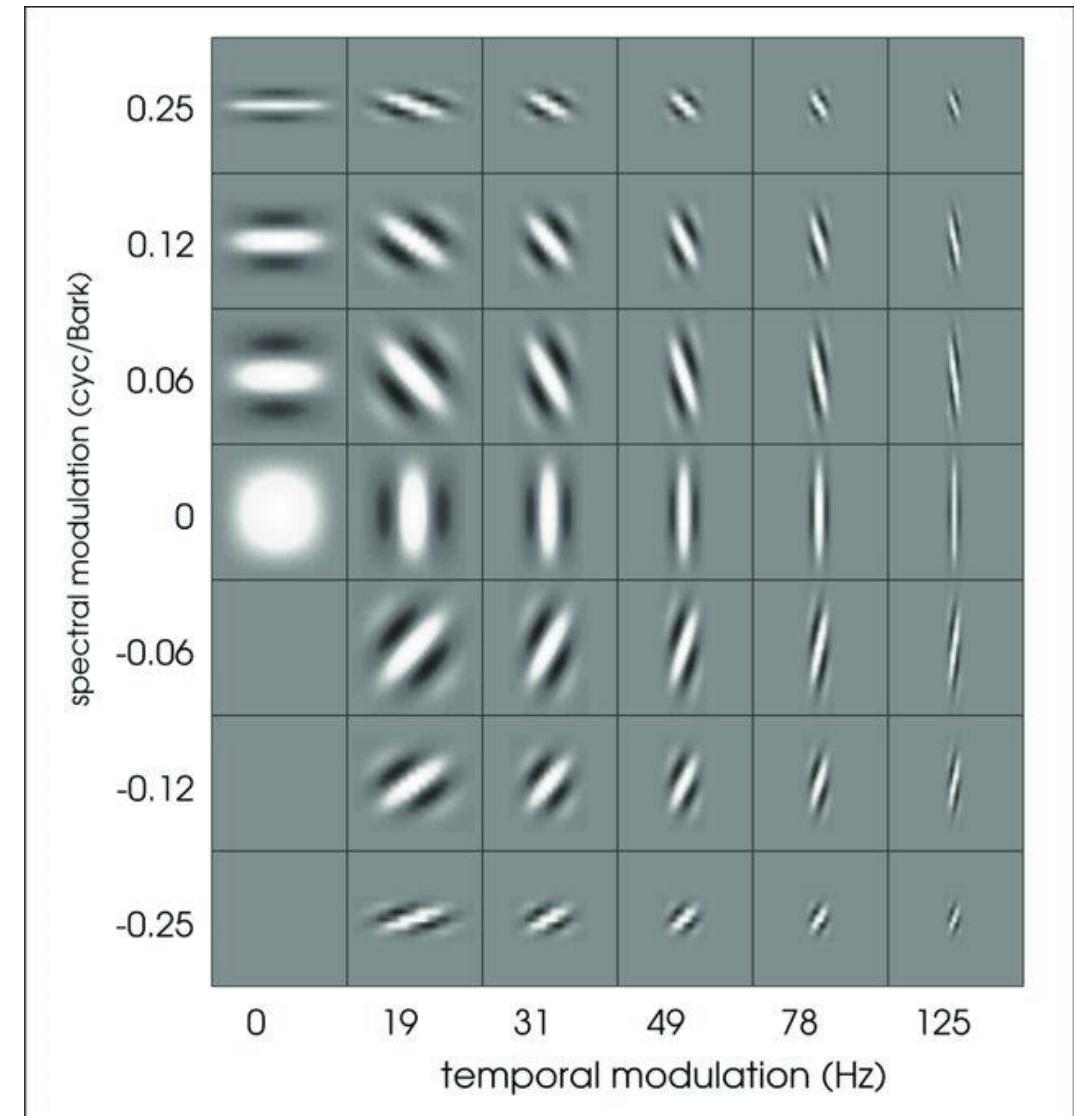
# Classical approaches

- Filtering (spatial domain)
- Frequency domain



Gabor filter

[https://www.researchgate.net/publication/228411540\\_Directional\\_Adaptive\\_WSSG\\_Filter/figures?lo=1](https://www.researchgate.net/publication/228411540_Directional_Adaptive_WSSG_Filter/figures?lo=1)

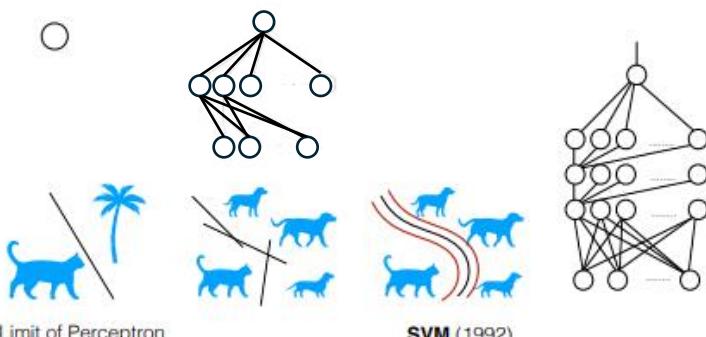


[https://www.researchgate.net/publication/313489901\\_Matching\\_Pursuit\\_Analysis\\_of\\_Auditory\\_Receptive\\_Fields%27\\_Spectro-Temporal\\_Properties/figures?lo=1](https://www.researchgate.net/publication/313489901_Matching_Pursuit_Analysis_of_Auditory_Receptive_Fields%27_Spectro-Temporal_Properties/figures?lo=1)

# Neural Network timeline

1 node  
multiple nodes  
shallow structure

multiple nodes  
deep structure



Limit of Perceptron (1969)

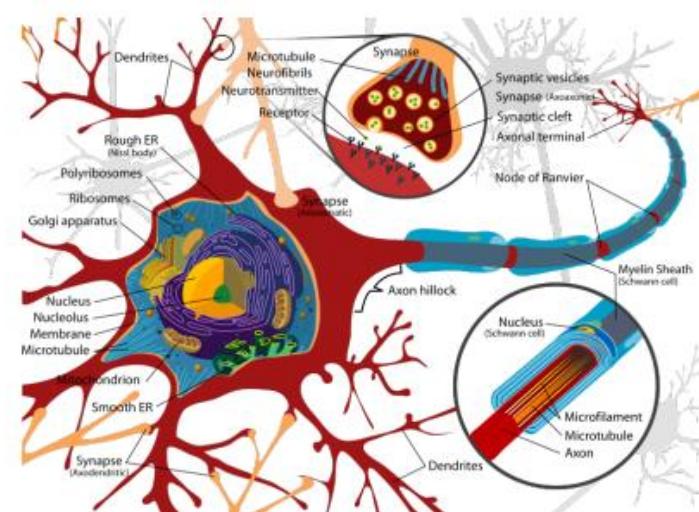
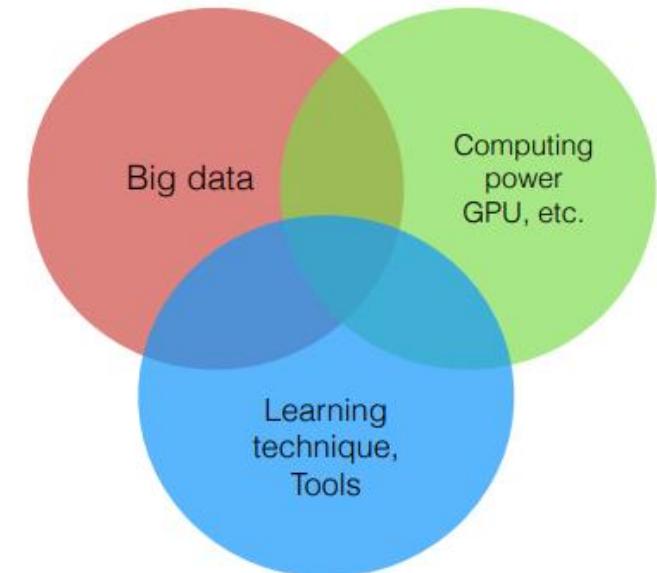
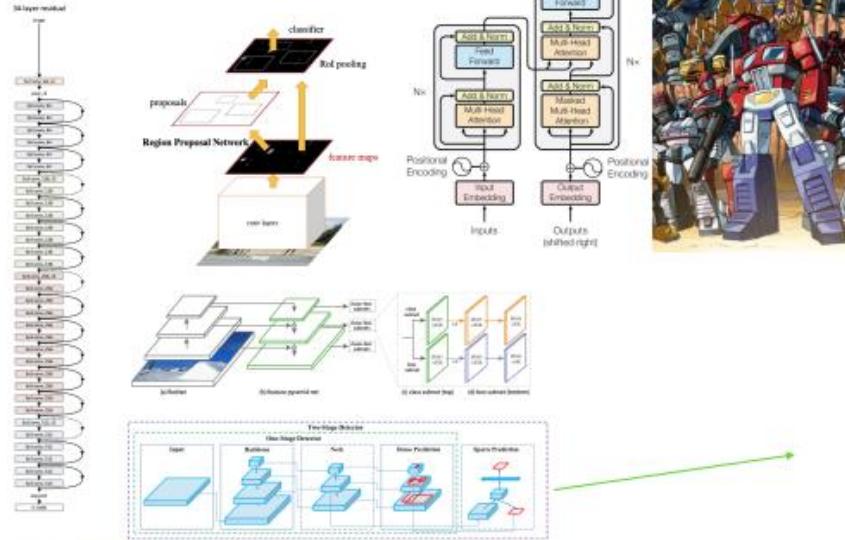
McCulloch & Pitts Neuron model (1940)

Rosenblatt's **Perceptron** (1957)

SVM (1992)

Backprop  
70's-80's

Deep Learning  
(2006)

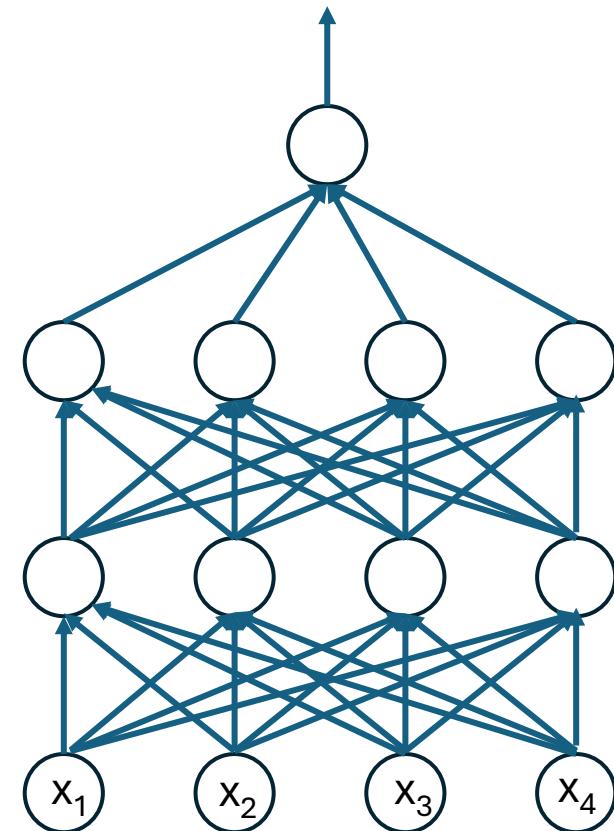


■ The Perceptron (Frank Rosenblatt at Cornell University, 1957)

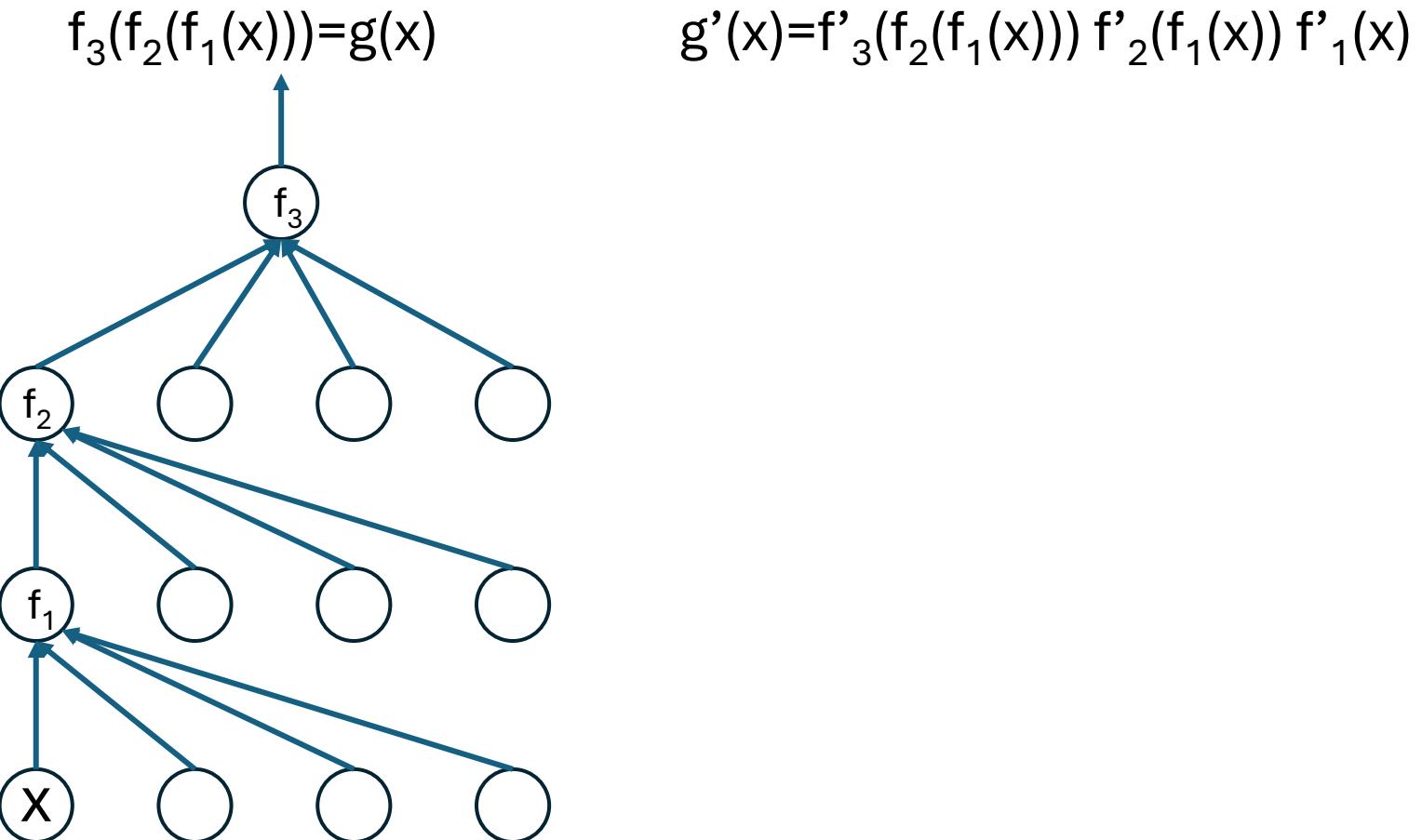


# Basic NN

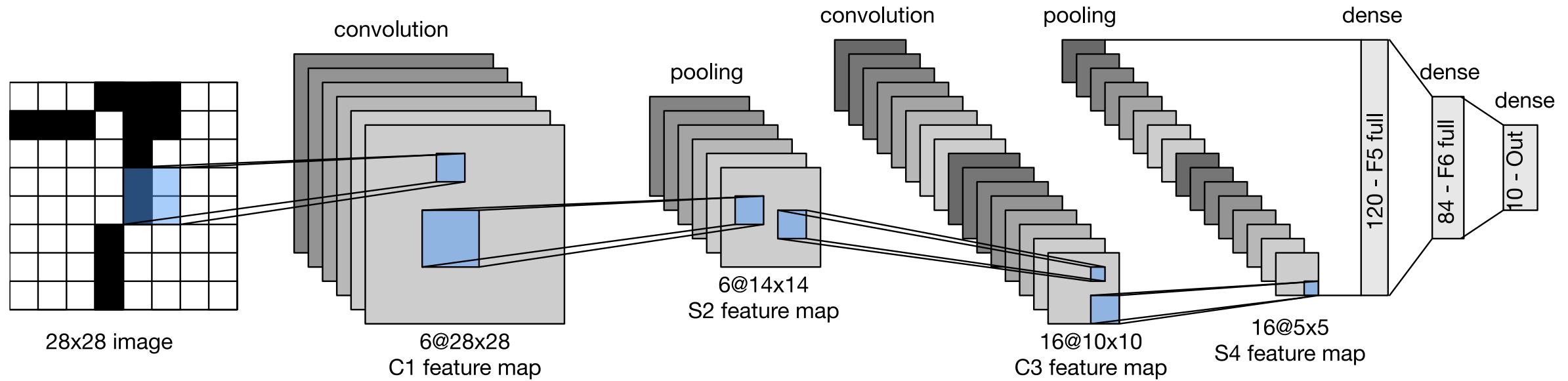
- Gradient of a function
  - Vector of partial derivatives of that function
  - Direction of changes that increases the value of the function
- Loss function
- Gradient descent
- Learning rate



# Gradient Vanishing problem



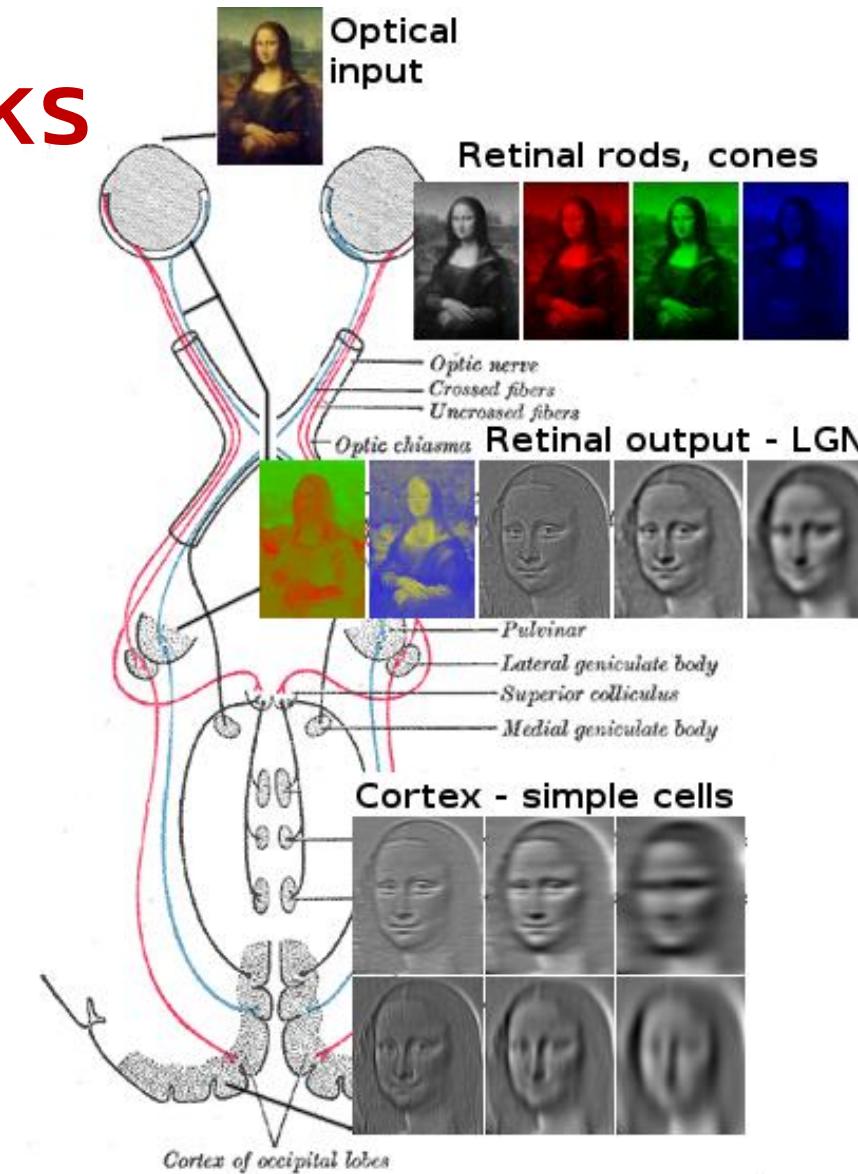
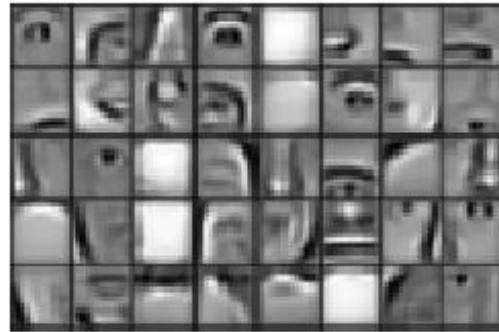
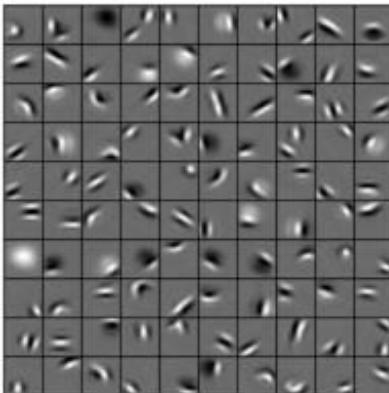
# Gradient vanishing solution 1: CNN



By Zhang, Aston and Lipton, Zachary C. and Li, Mu and Smola, Alexander J. -  
<https://github.com/d2l-ai/d2l-en>, CC BY-SA 4.0,  
<https://commons.wikimedia.org/w/index.php?curid=152265656>

# Convolutional Neural Networks

- Classical machine learning approach is composed of
  - Feature extractor
  - Classifier
- Deep learning extracts features
  - Less heuristic
  - Hierarchical representation

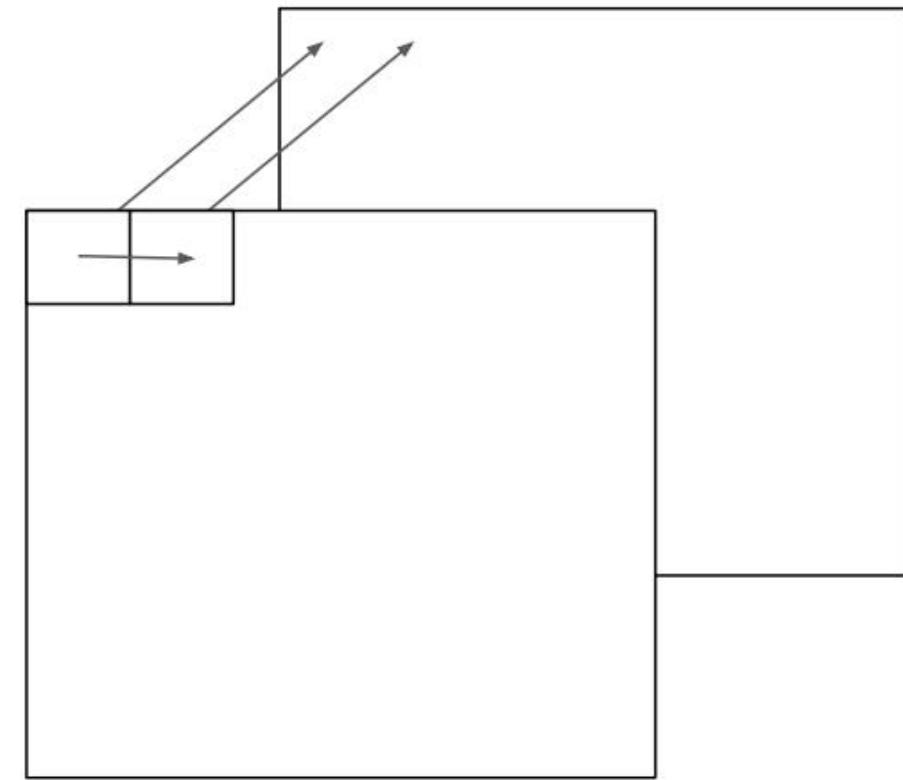




**Input image**

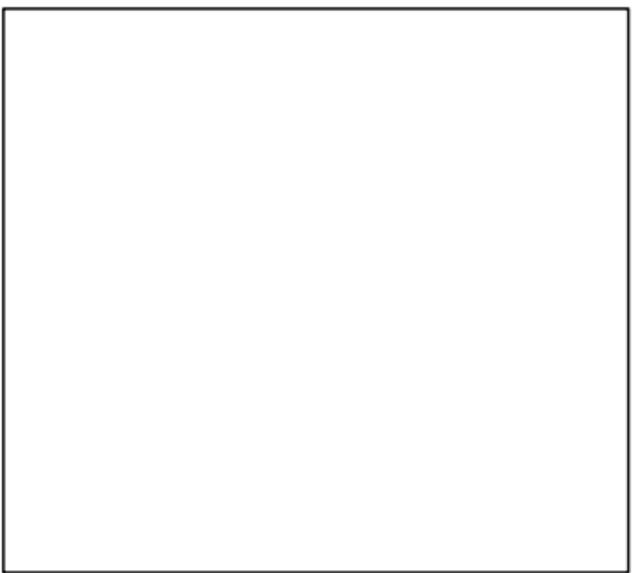


**Convolution Filter**



**Input image**

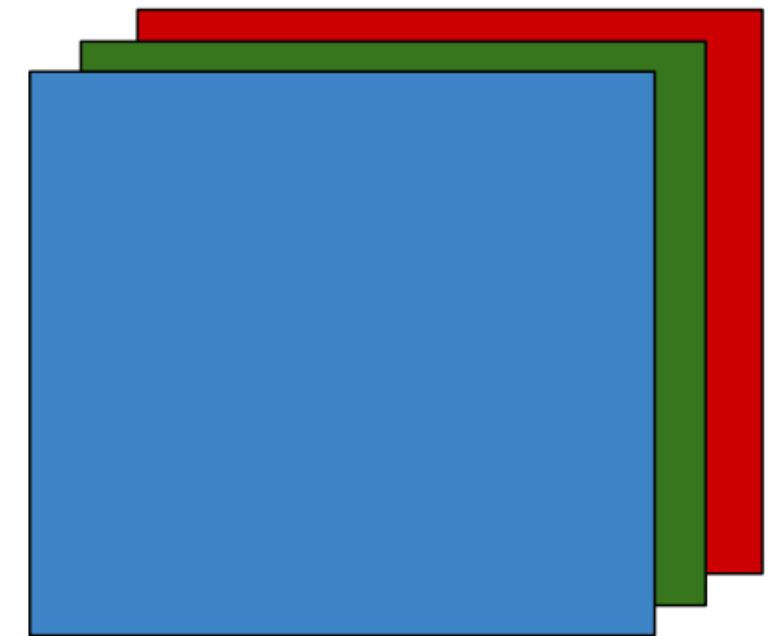
**Feature map**



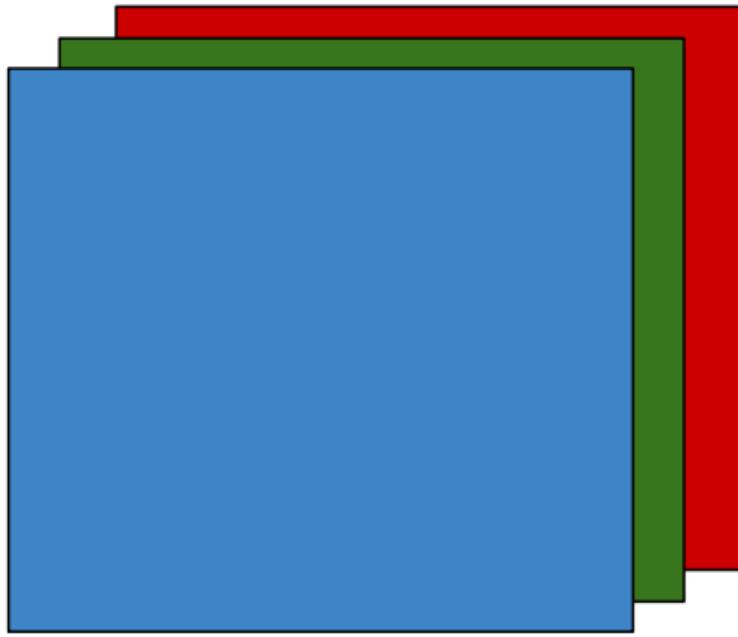
**Input image**



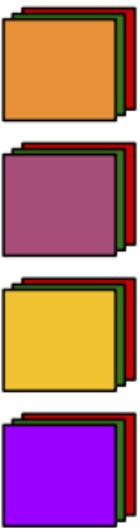
**Convolution Filter**



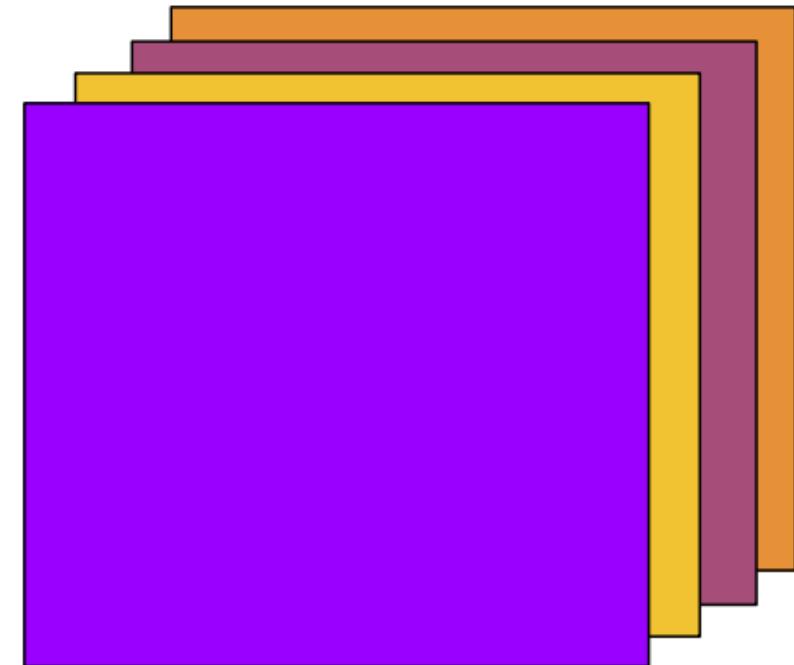
**Feature map**



Input feature maps



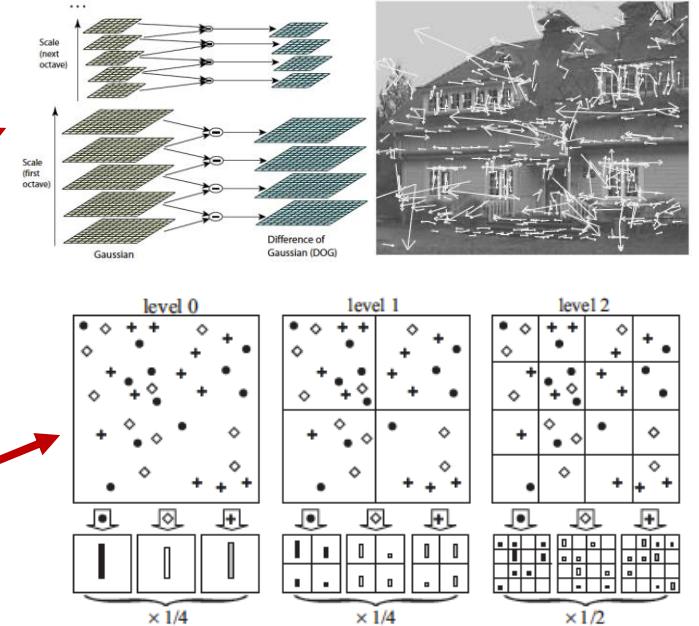
Convolution Filter



Output feature maps

- In the past, convolution filters needed to be prepared first
- Now, convolution filters = Weight of a neural network
- Weight sharing strategy
  - Need to sum all gradient before update
    - Why this helps gradient vanishing?

# Deep Features



- GIST (global feature) + SVM (RBF): 85.57%
- SIFT (local feature) + BoF + SVM (Histogram intersection): 89.69%
- SIFT + SPM (spatial pyramid matching) + LLC (locality-constrained linear coding) + SVM (linear): 91.48%
- CNN (AlexNet trained on other dataset) + SVM (linear): 93.58%

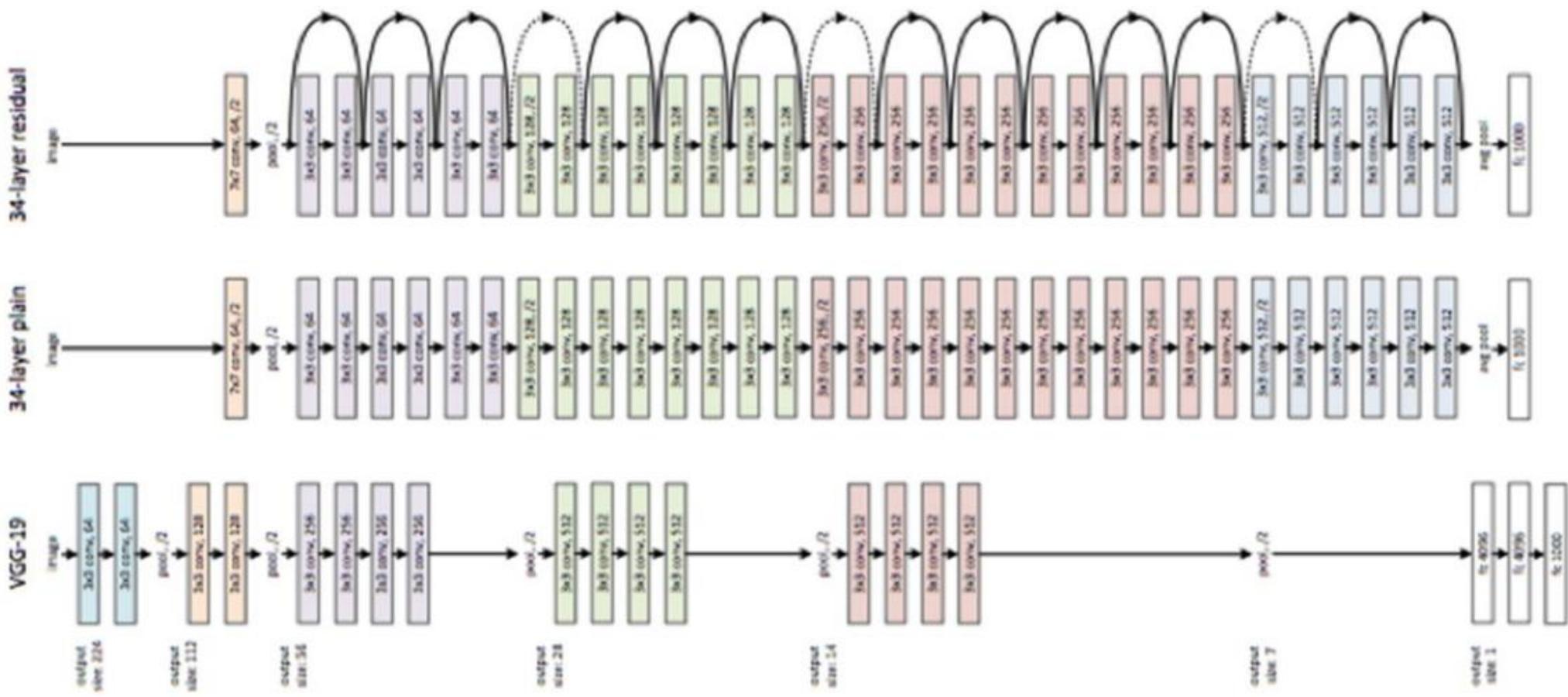
$$\min_c \sum_{i=1}^M \|x_i - Bc_i\|^2 + \lambda \|d_i \odot c_i\|^2$$

subject to  $1^T c_i = 1$

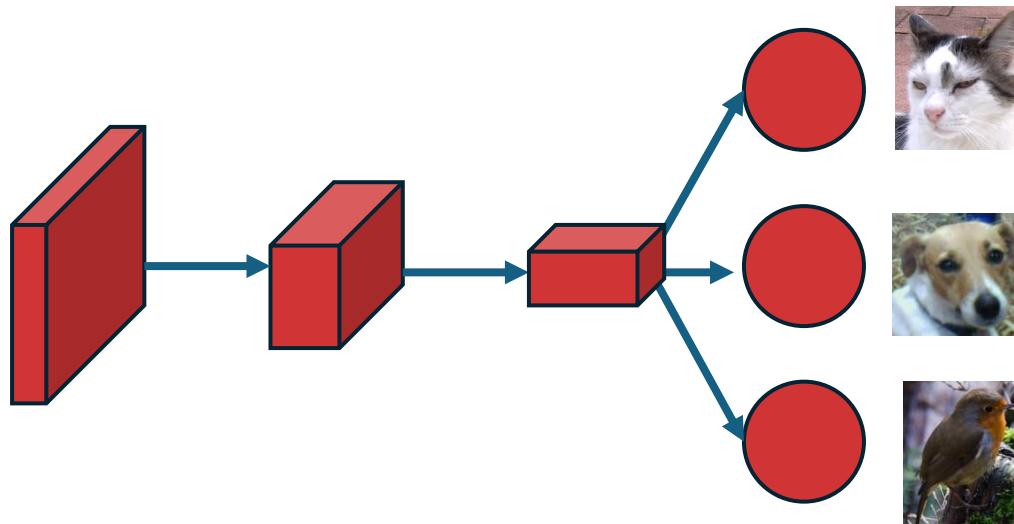
$$d_{im} = \exp\left(\frac{1}{Z\sigma}\|x_i - b_m\|^2\right)$$
$$Z = \max_k \|x_i - b_k\|^2,$$

# Gradient Vanishing solution 2: ResNet / skip connection

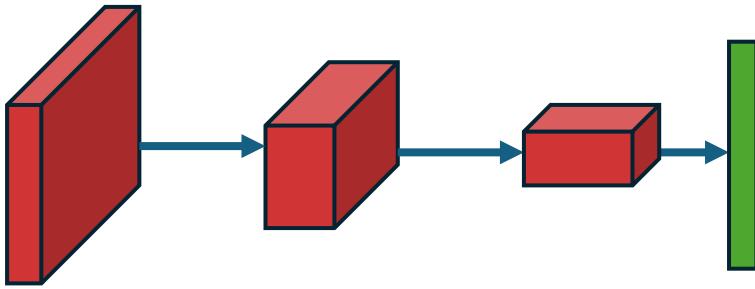
skip connection



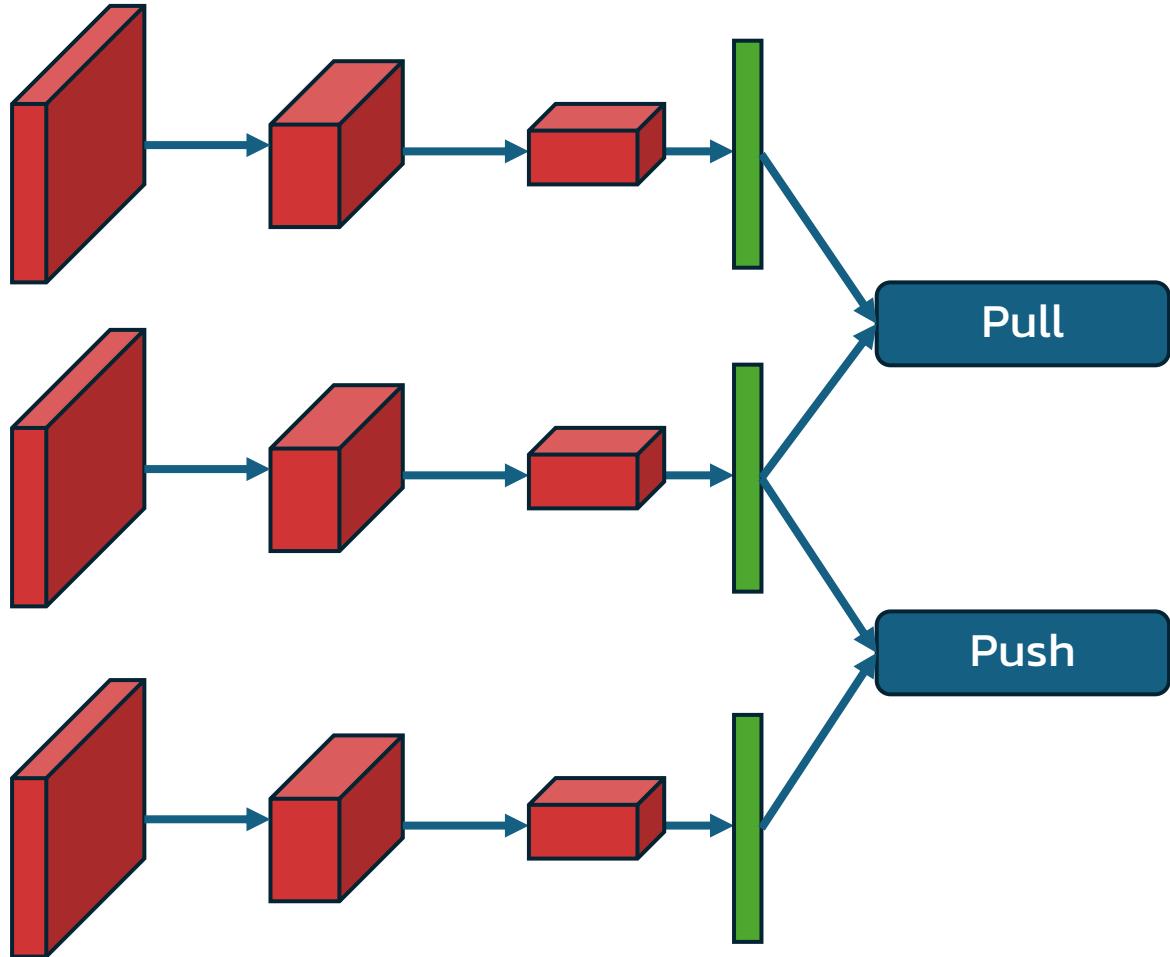
# **Some common architectures & tasks**



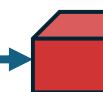
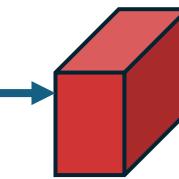
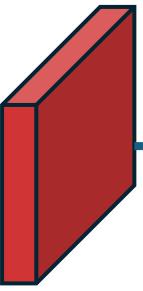
Classical CNN for Image Classification  
Chest X-ray, Thalassemia



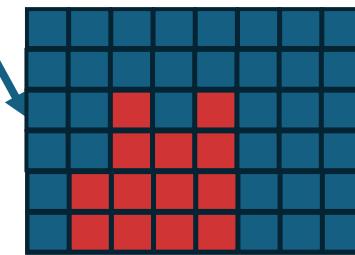
Classical CNN for Feature Extraction



# Feature for Face Verification



Bounding box + Confidence

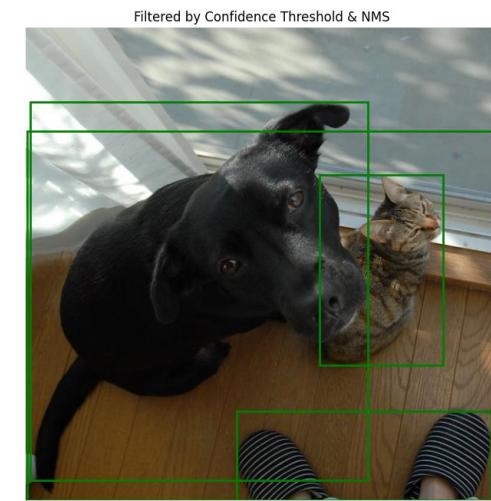
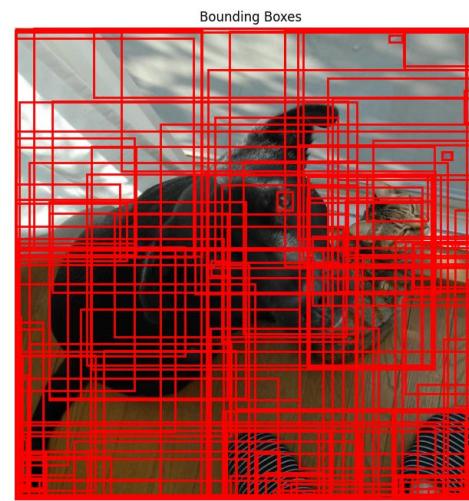
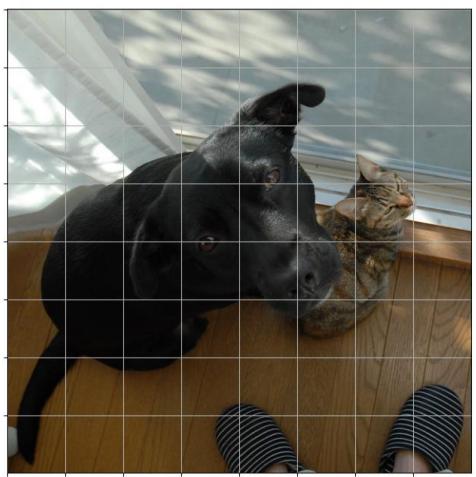


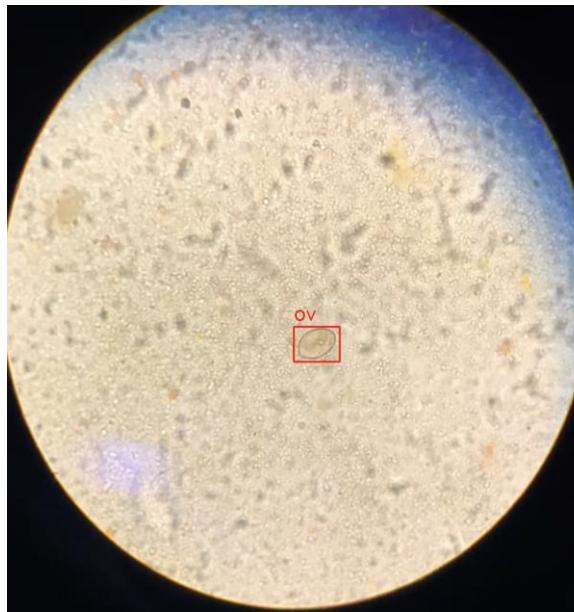
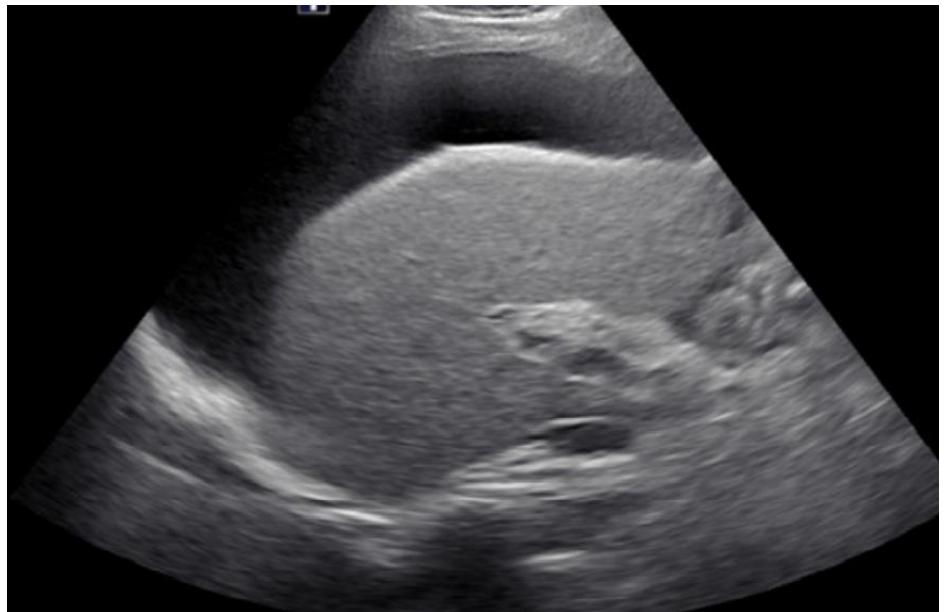
Class posterior probability

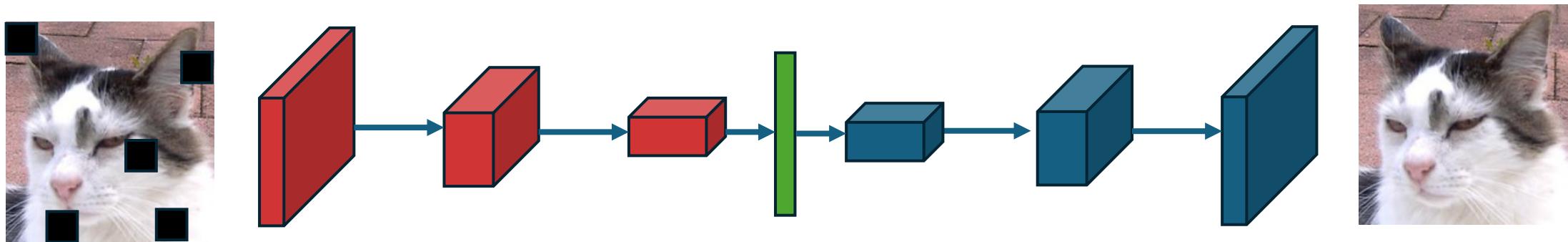


Detection result

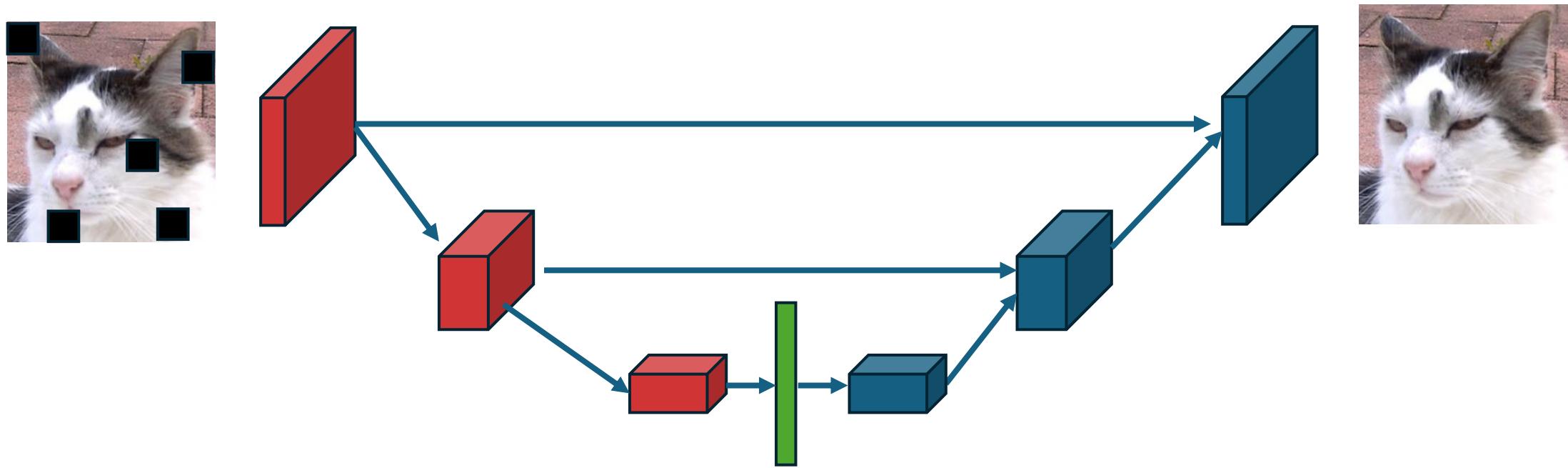
YOLO for Object Detection



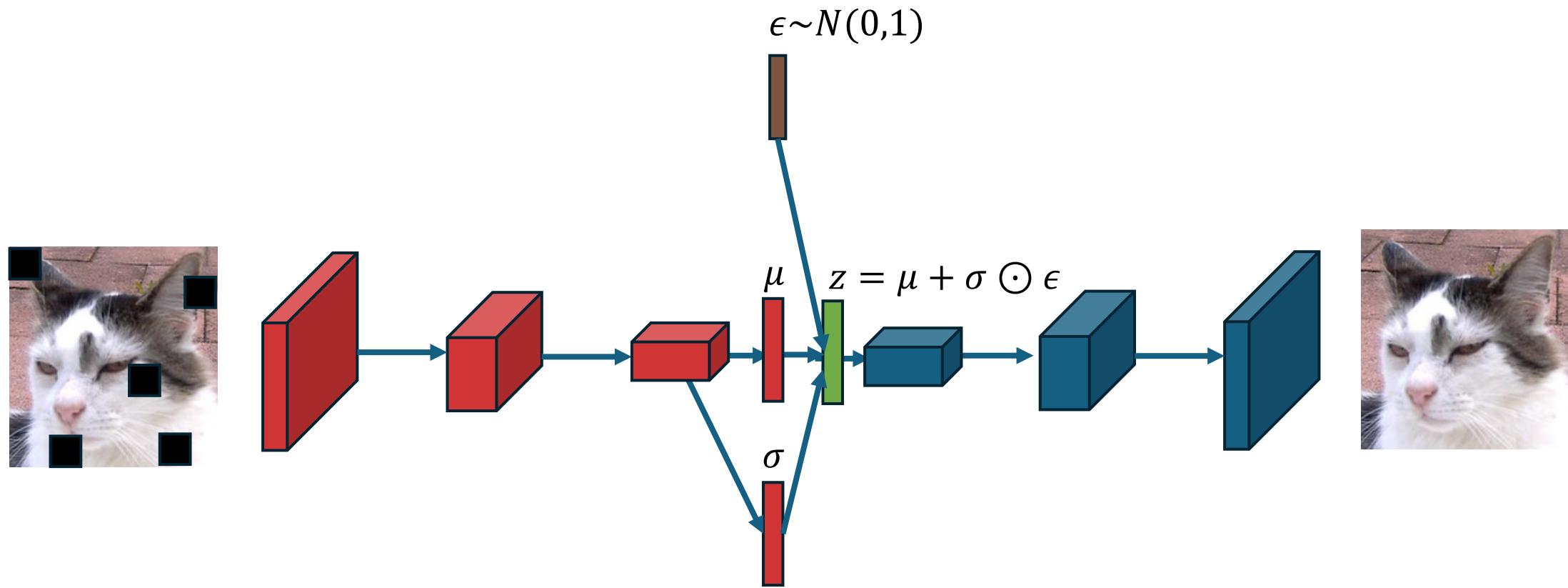




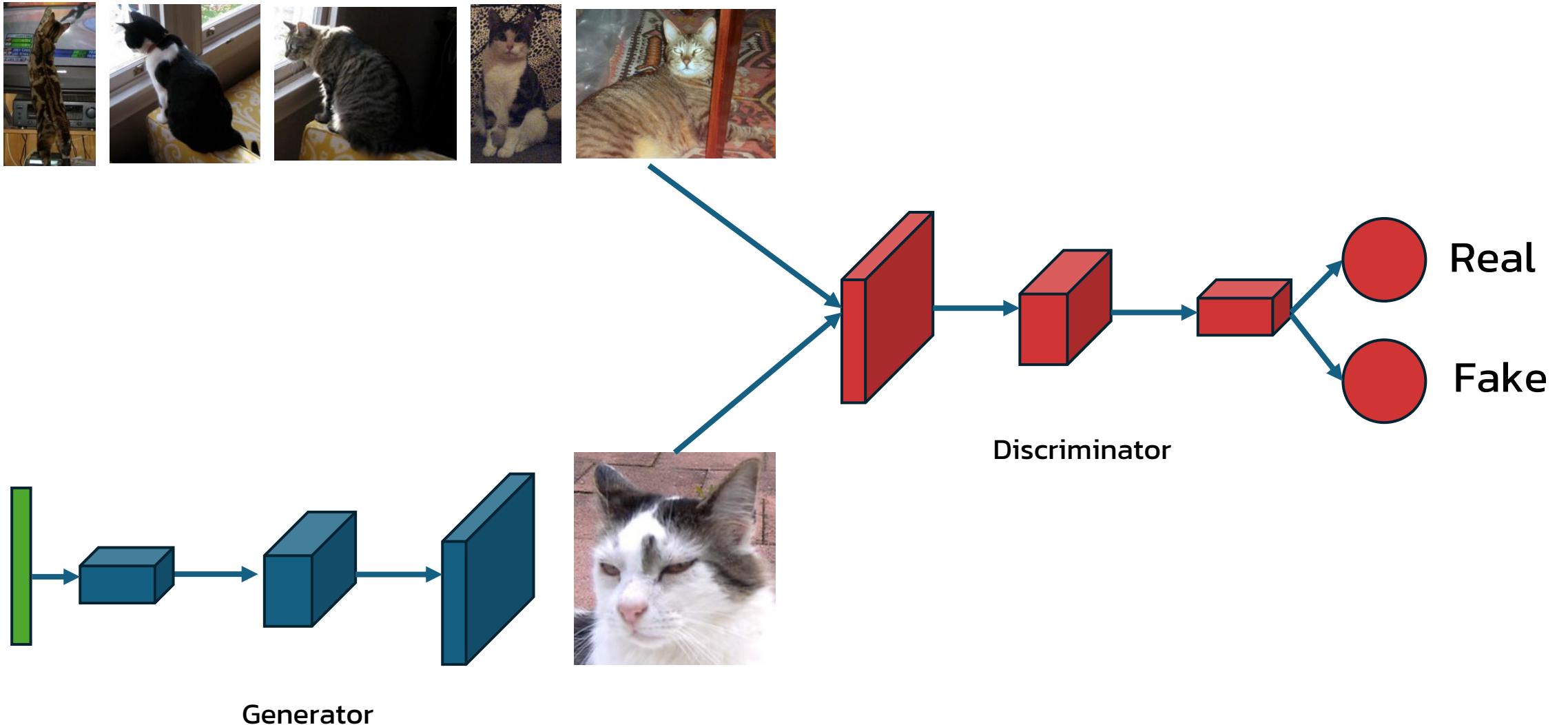
Encoder-Decoder for  
Image Denoising, Image Segmentation



U-Net for  
Image Denoising, Image Segmentation, Diffusion Model

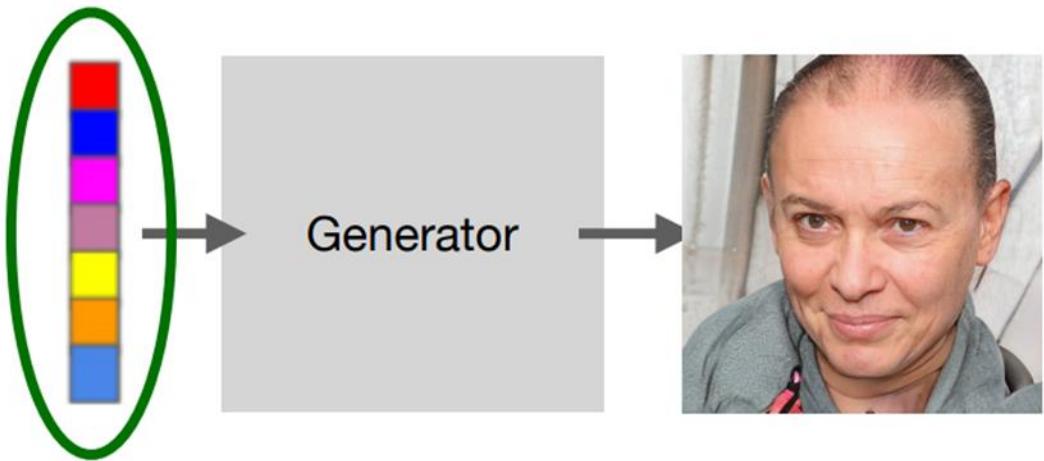


Variational Autoencoder (VAE) for  
Image Generation

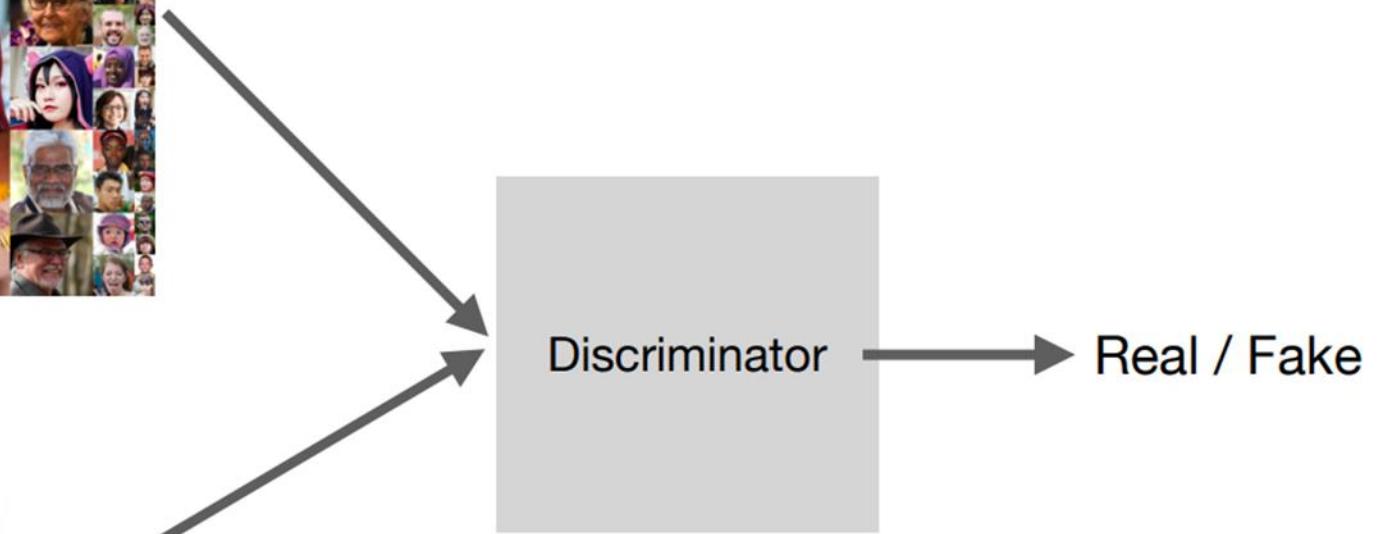


**Generative Adversarial Networks (GAN) for  
Image Generation, Super Resolution**

# GAN



Adjust face appearance





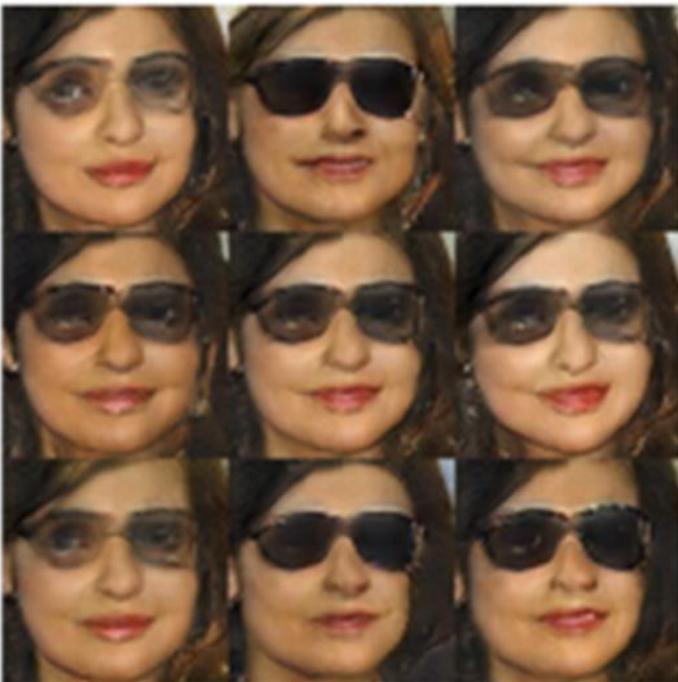
man  
with glasses



man  
without glasses



woman  
without glasses



woman with glasses

A. Radford et al. "Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks", 2016

# X does not exist



This sneaker does not exist

<https://thissneakerdoesnotexist.com/>

ENTIRE GUEST SUITE  
**Unique Finca. Beautiful, calm nest!**

Paris

4 guests 2 bedrooms 2 beds 2 baths

Private bachelor home designed in modern, has all amenities. When it is cozy and clean fully furnished studio age of 2013. And the shaped romantic aristos hands to suit your stay... all amenities. Soft laundr

Anais

A screenshot of a website for a non-existent sneaker store. The header features a navigation bar with links for "The Grid", "Sneaker Editor", "Info", and "Contact", along with social media icons for Instagram and Twitter. Below the header, there are four pairs of different sneakers displayed horizontally. The main content area contains the URL "https://thissneakerdoesnotexist.com/" followed by a grid of five smaller images showing various interior rooms of a house, including a bedroom, a kitchen, and a living room with a fireplace.

<https://thispersondoestnotexist.com/> <https://thiscatdoesnotexist.com/>

<https://thisrentaldoesnotexist.com/>

**Another gradient vanishing (for sequence)**

# Classical Language Model

- “The wheel on the bus goes ...”
- **N-gram model**  $p(x_t|x_{<t}) = p(x_t|x_{t-(n-1)}, \dots, x_{t-1})$ 
  - Each word depends on n-1 previous words
  - Unigram, Bigram, Trigram
  - Probability table, discrete observations
- Improved models: Hidden Markov Model (HMM), Conditional Random Field (CRF)
  - Still rely on discrete observations

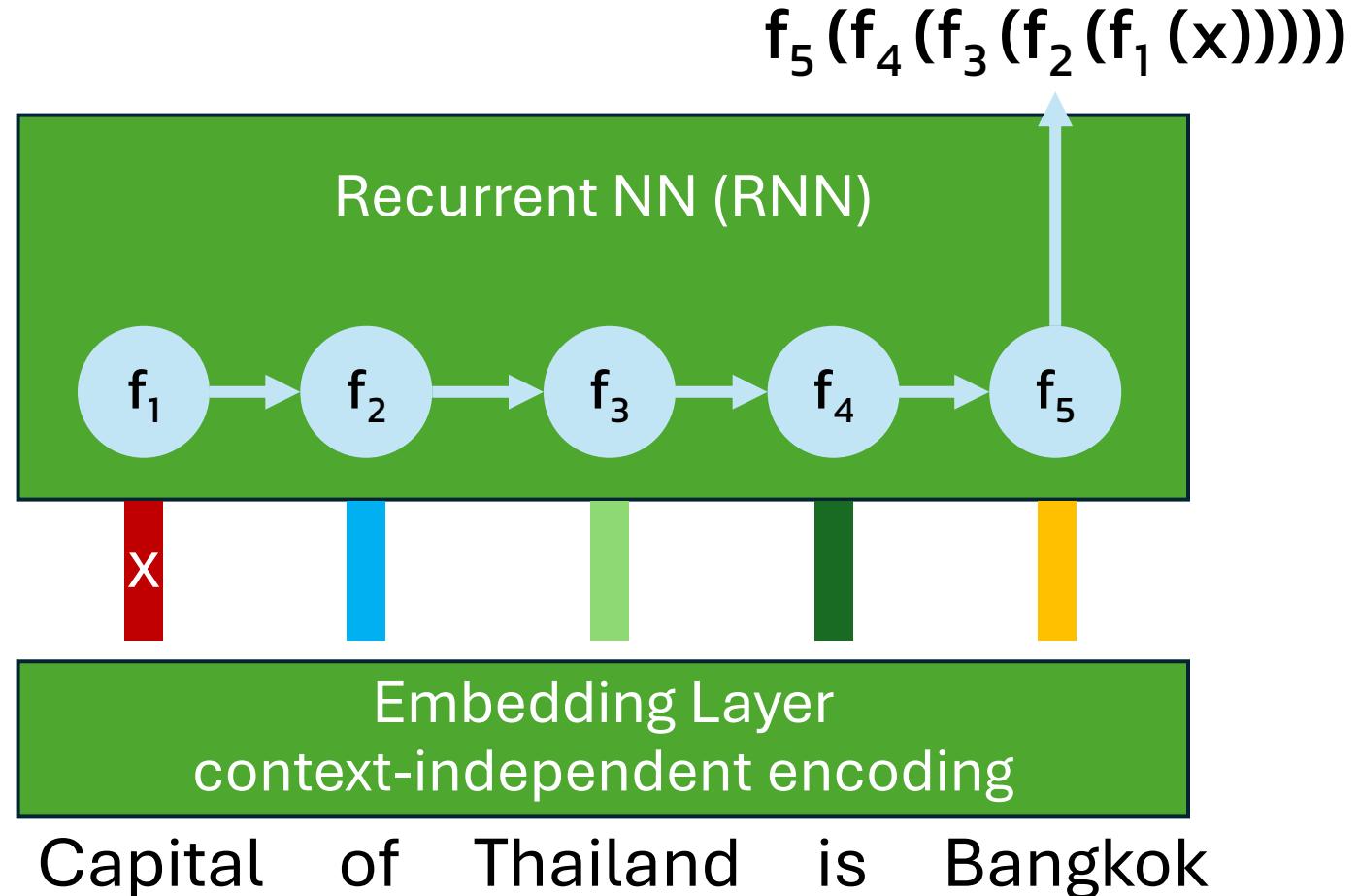
# Continuous representation

- **One-hot vector**
  - Dimension = number of words in vocabulary
  - 1 on position of this word in the vocab, 0 on every other dimensions
- **Bag-of-words** = present/absent of words in the input document
- **Latent Semantic Analysis (LSA)**= **Principal Component Analysis (PCA)** on bag-of-words vectors
  - Obtain a dense vector for each word
- **Word2vec**
  - Autoencoder for one-hot vector, with 1 encoding layer & 1 decoding layer
  - Train either to predict surrounding words or to use surrounding words to predict this word
  - After training, the latent from encoding layer = vector represent the word
  - $\text{vec}(\text{King}) - \text{vec}(\text{Man}) + \text{vec}(\text{Woman}) = \text{vec}(\text{Queen})$

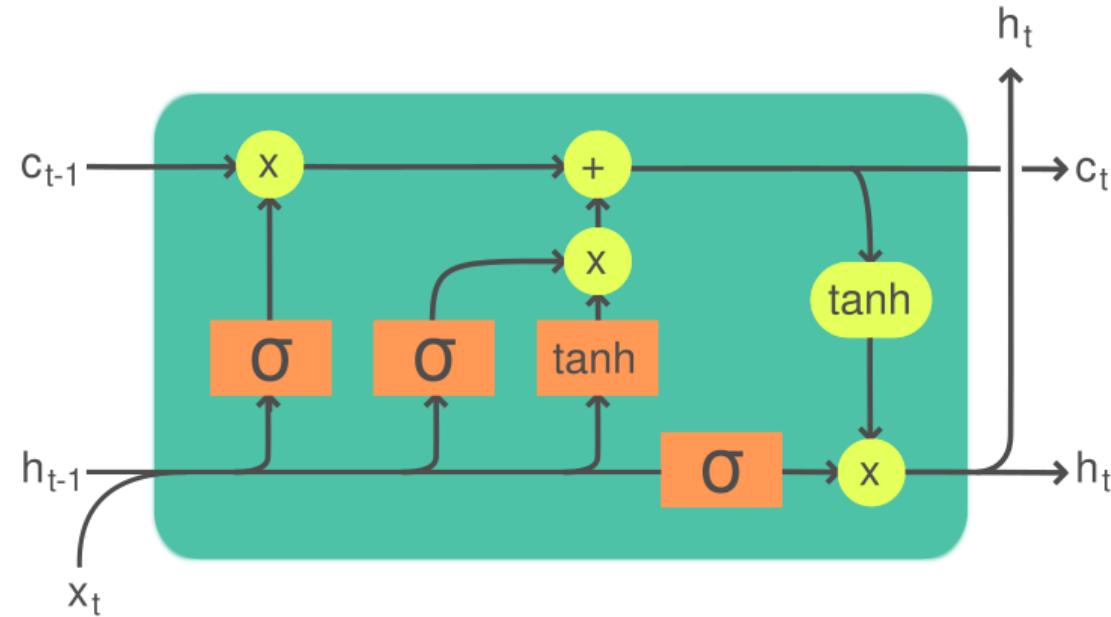
# Context-independent representation



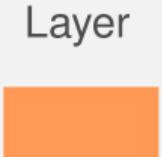
# Context-dependent representation



# Long Short-Term Memory (LSTM)



Legend:



Layer

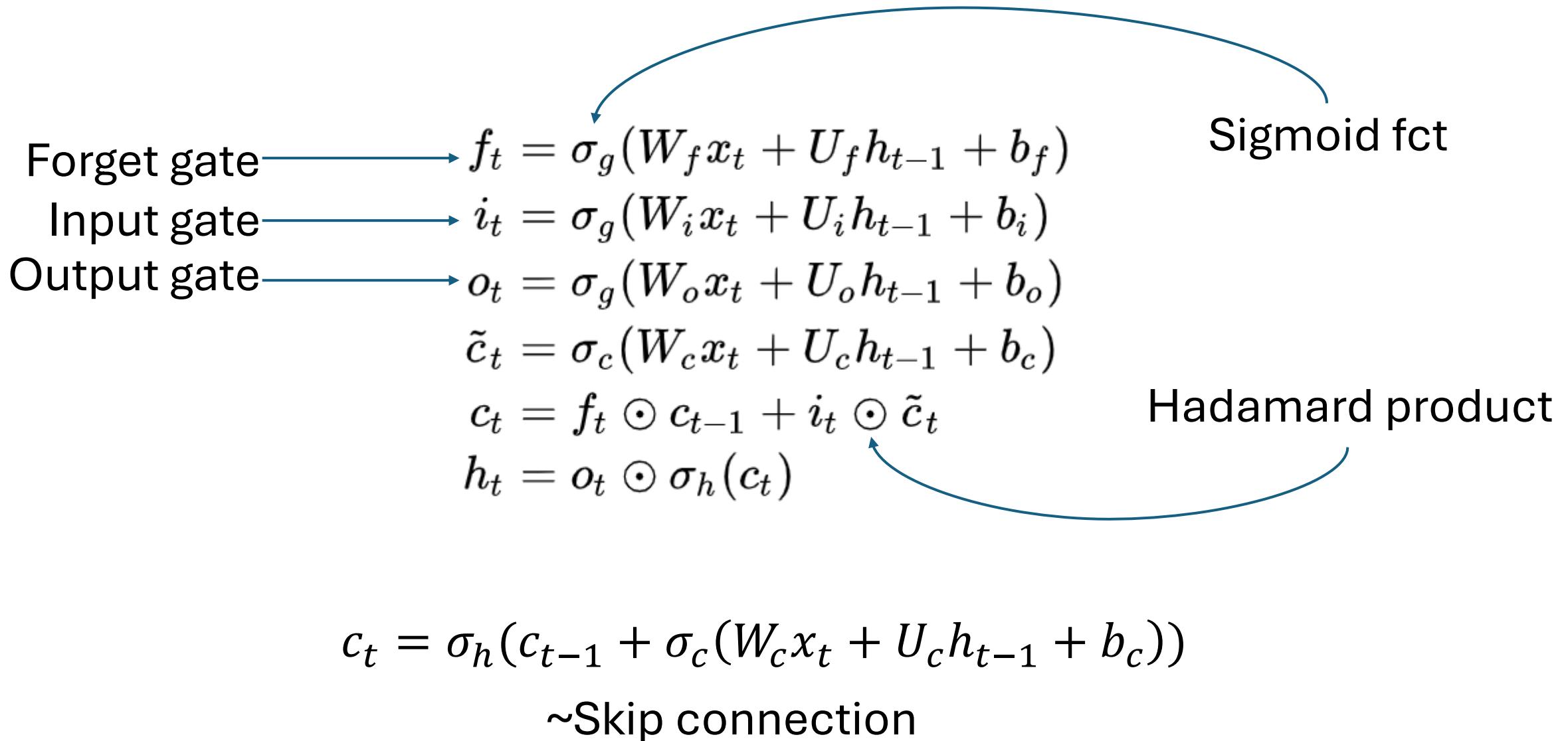
Componentwise

Copy

Concatenate

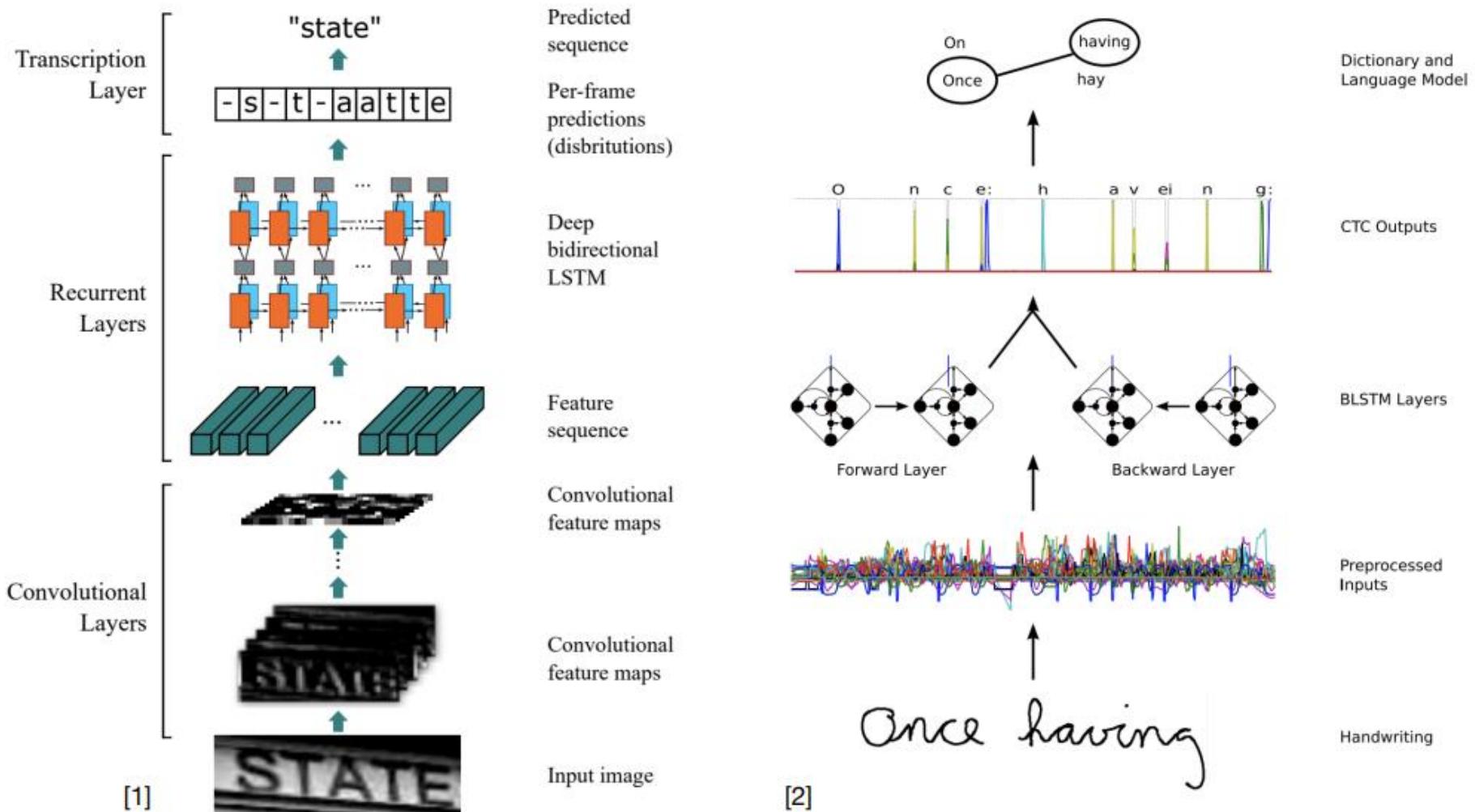


# LSTM





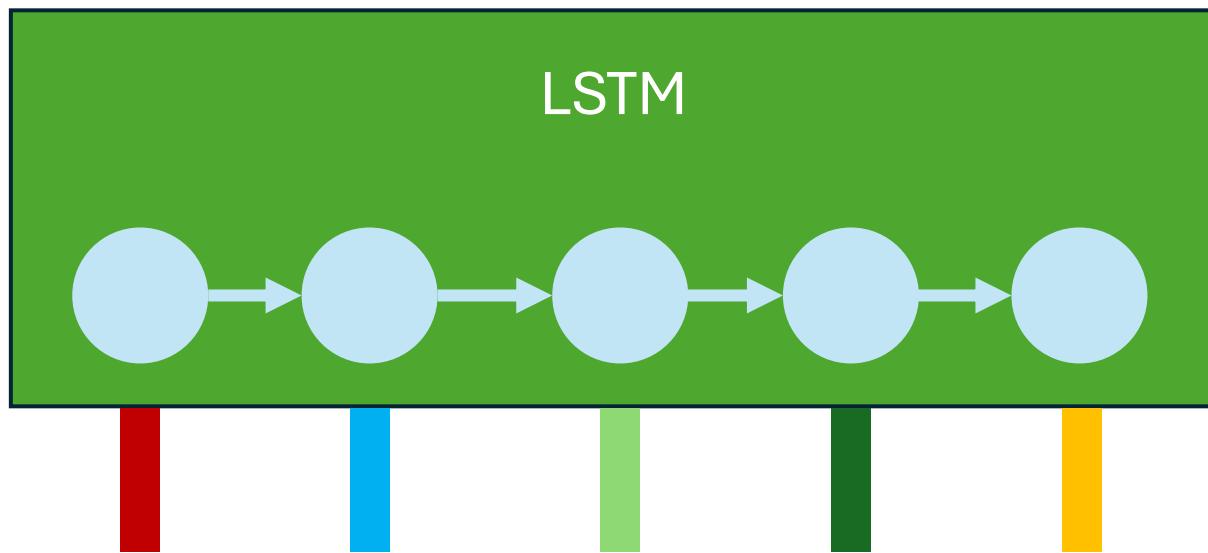
W. Byeon et al. "Scene Labeling with LSTM Recurrent Neural Networks", CVPR 2015



[1] B. Shi, X. Bai, and C. Yao, "An End-to-End Trainable Neural Network for Image-based Sequence Recognition and Its Application to Scene Text Recognition", <https://arxiv.org/pdf/1507.05717.pdf>

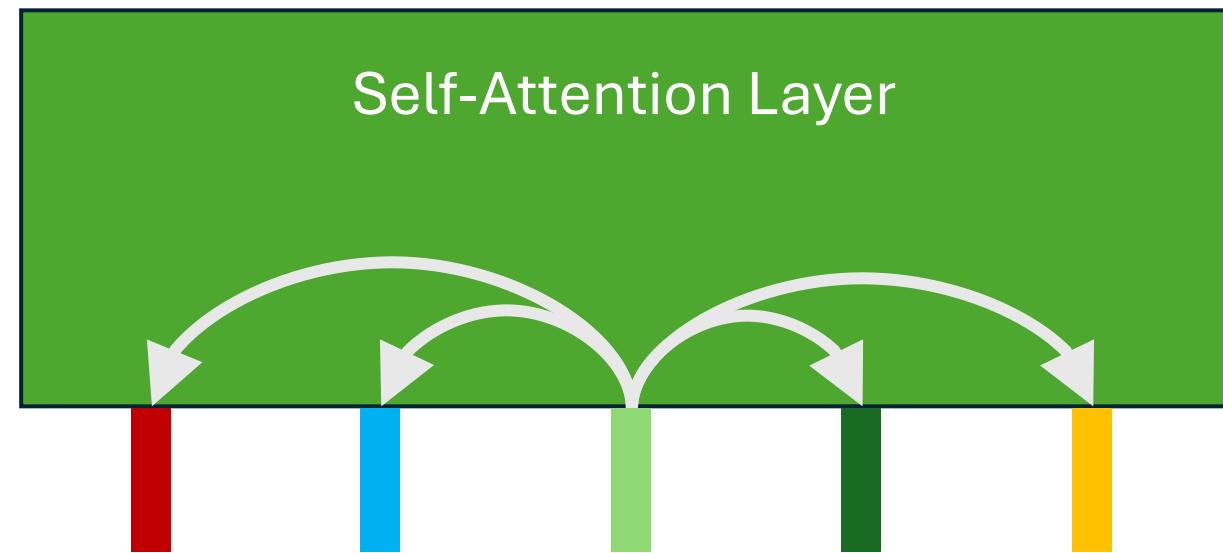
[2] A. Graves, M. Liwicki, S. Fernández, R. Bertolami, H. Bunke, and J. Schmidhuber, "A Novel Connectionist System for Unconstrained Handwriting Recognition", [http://www.cs.toronto.edu/~graves/tpami\\_2009.pdf](http://www.cs.toronto.edu/~graves/tpami_2009.pdf)

# Context-dependent representation



Embedding Layer  
context-independent encoding

Capital of Thailand is Bangkok

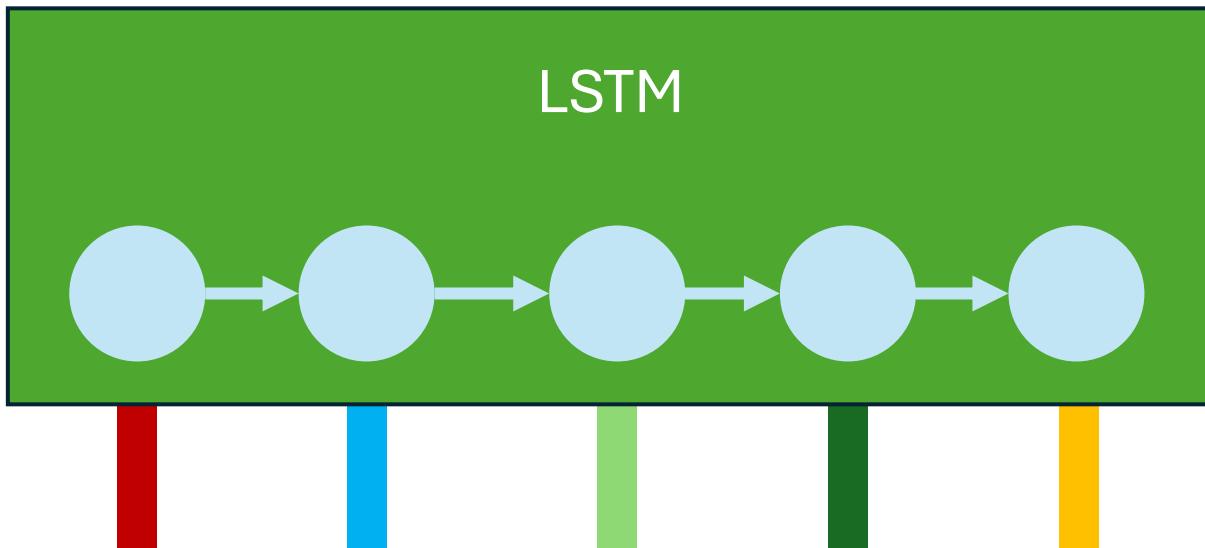


Embedding Layer  
context-independent encoding

Capital of Thailand is Bangkok

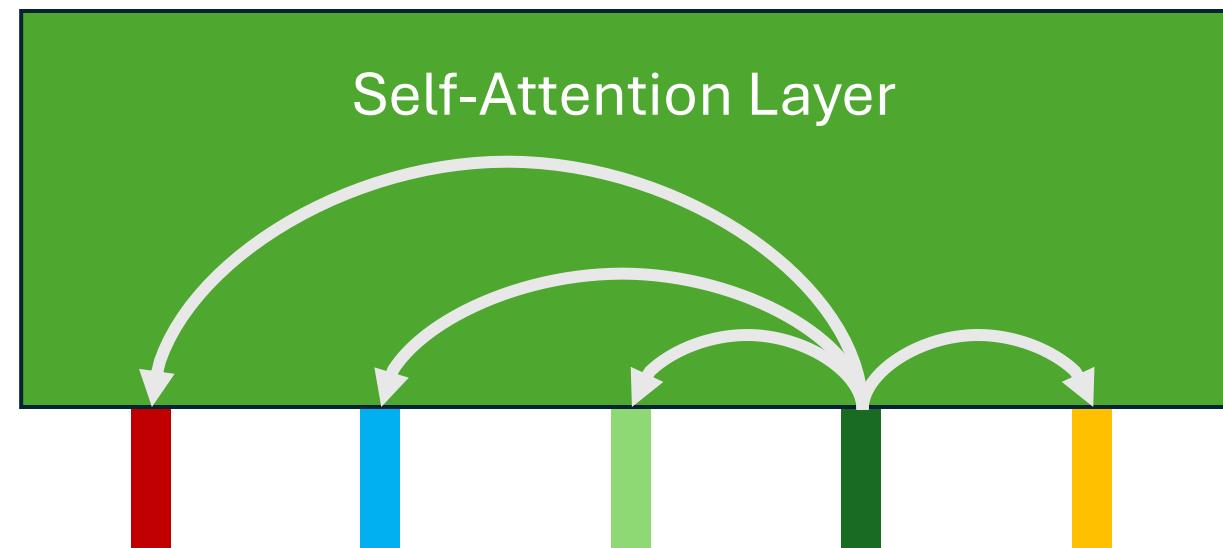
# Gradient vanishing problem

- Self-attention mechanism solves gradient vanishing with increasing computational cost



Embedding Layer  
context-independent encoding

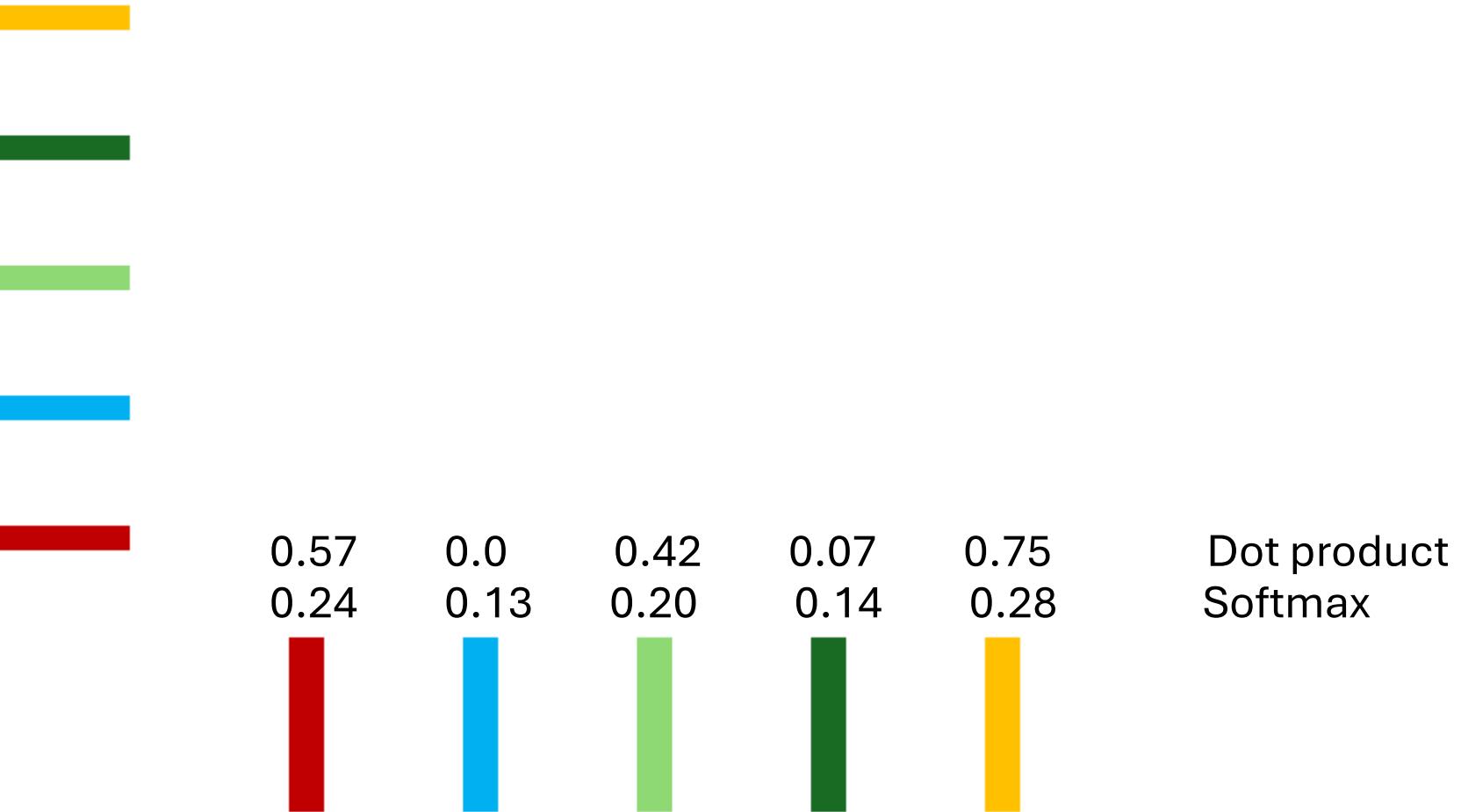
Capital of Thailand is Bangkok



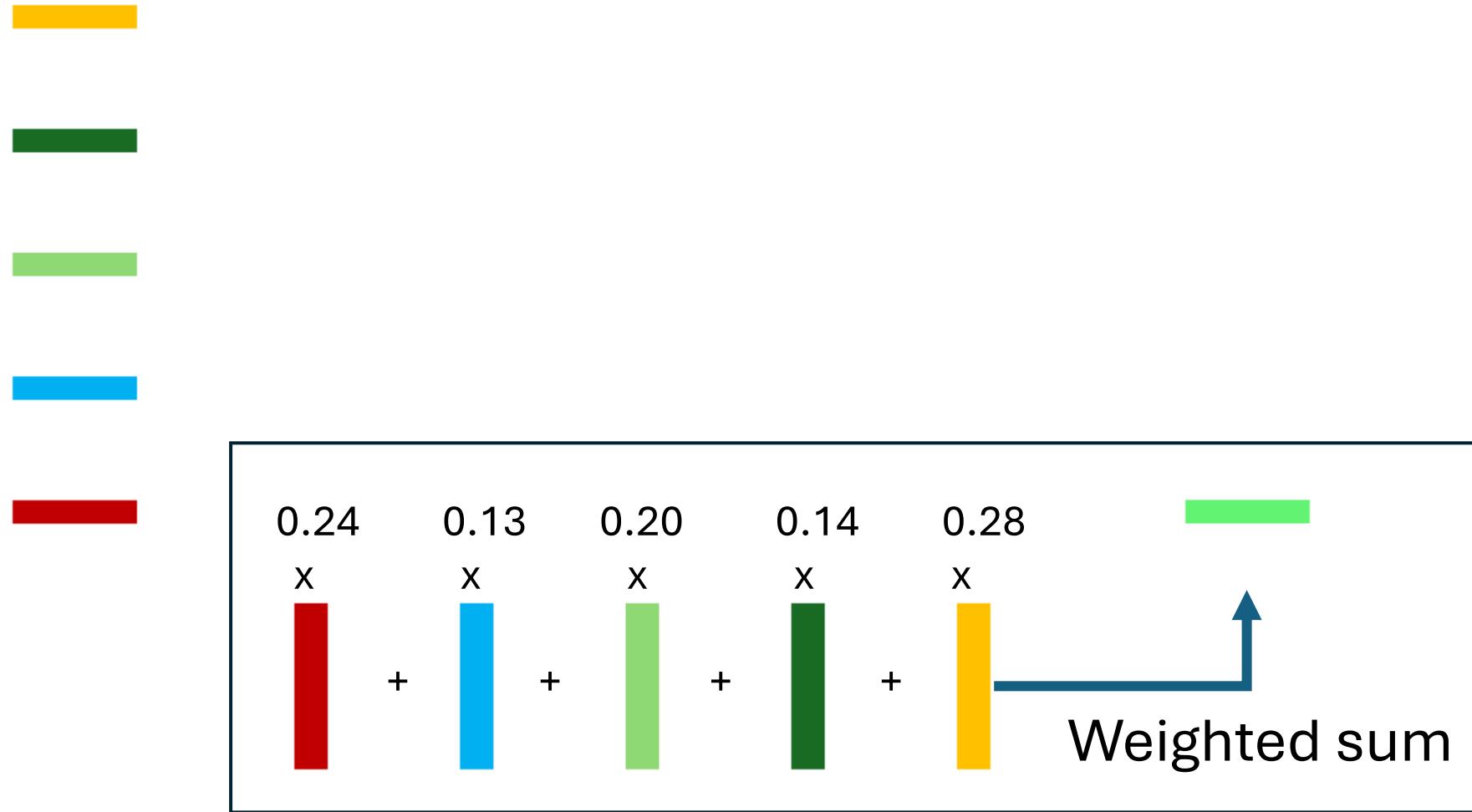
Embedding Layer  
context-independent encoding

Capital of Thailand is Bangkok

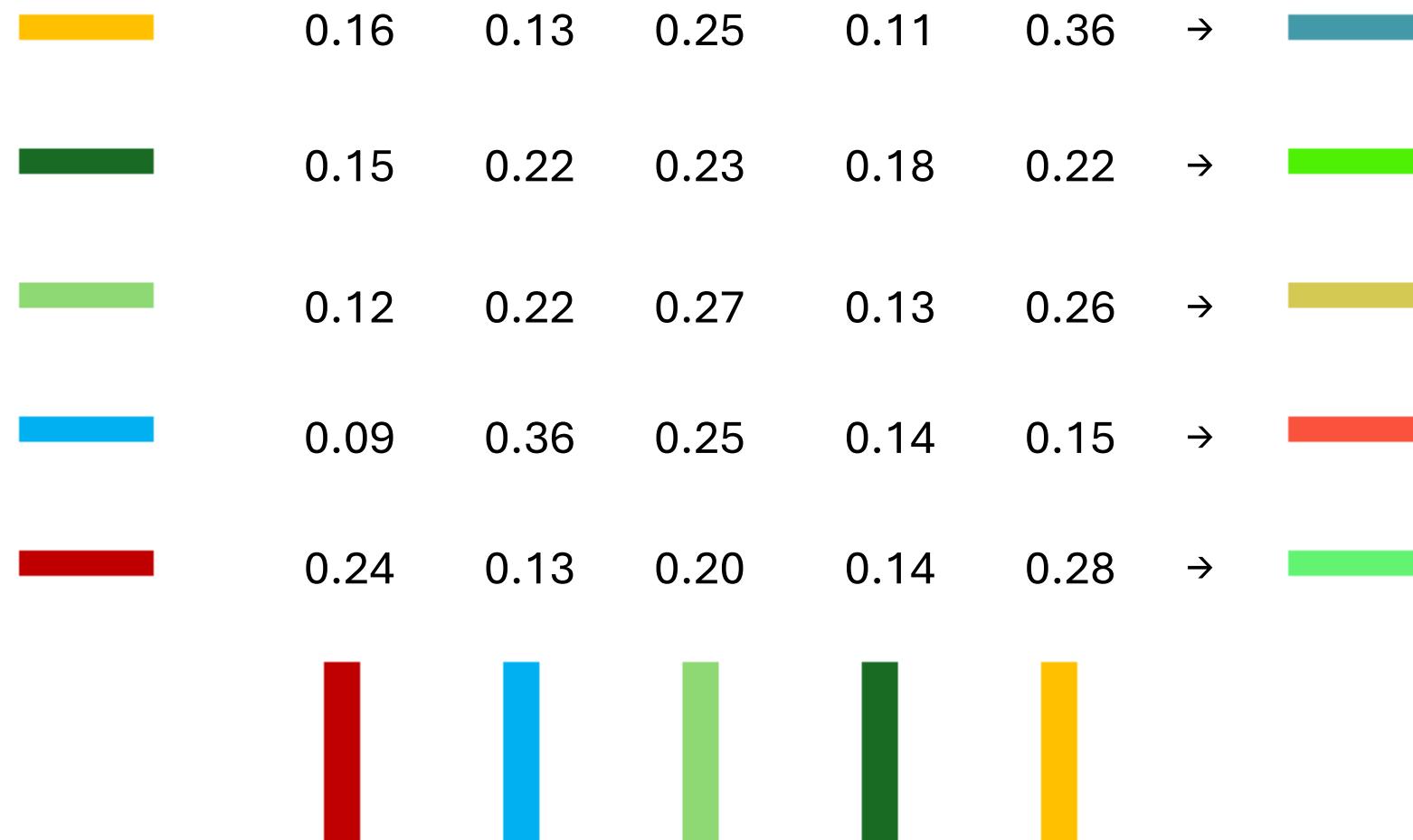
# Self-Attention Layer



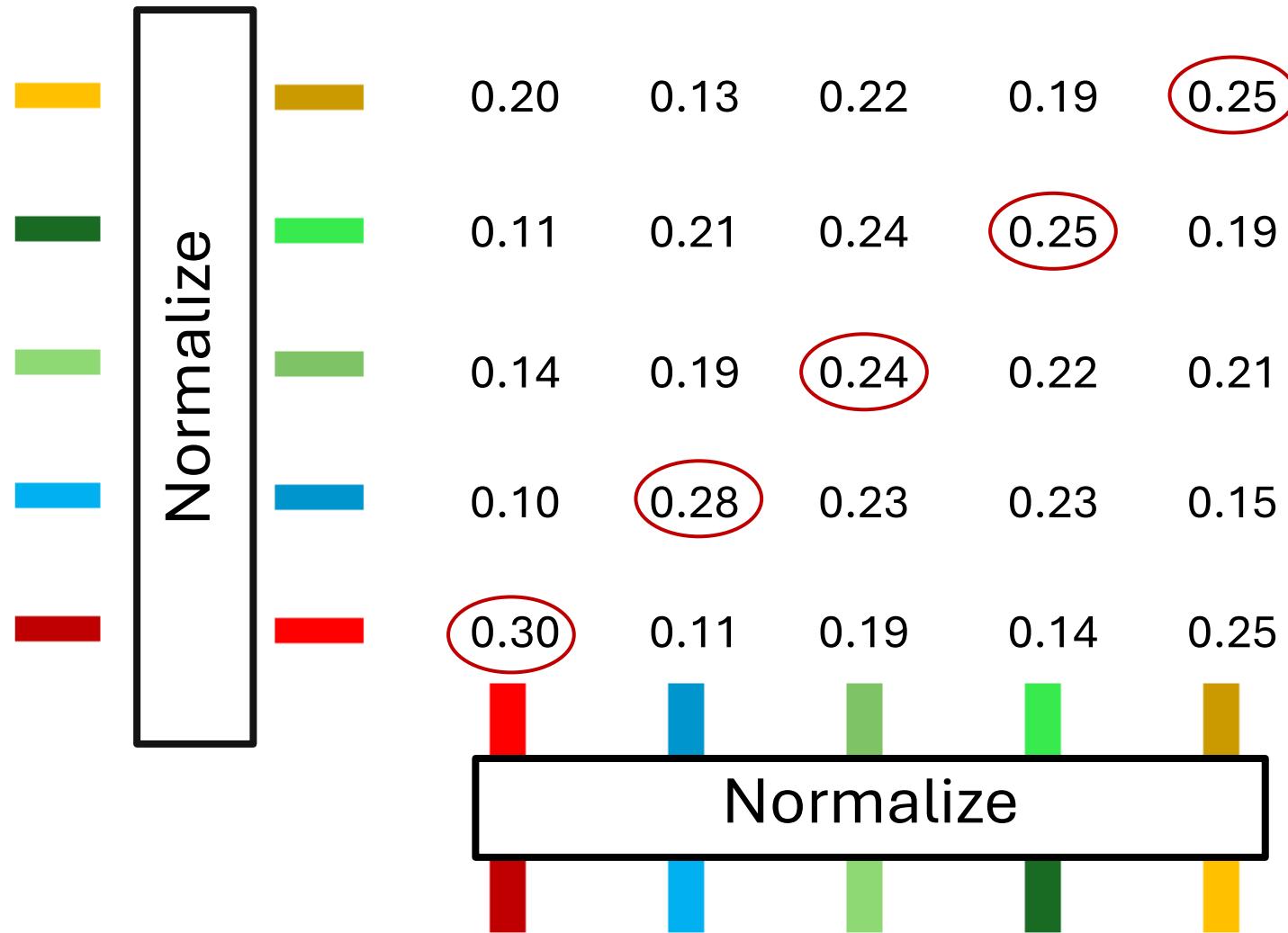
# Self-Attention Layer



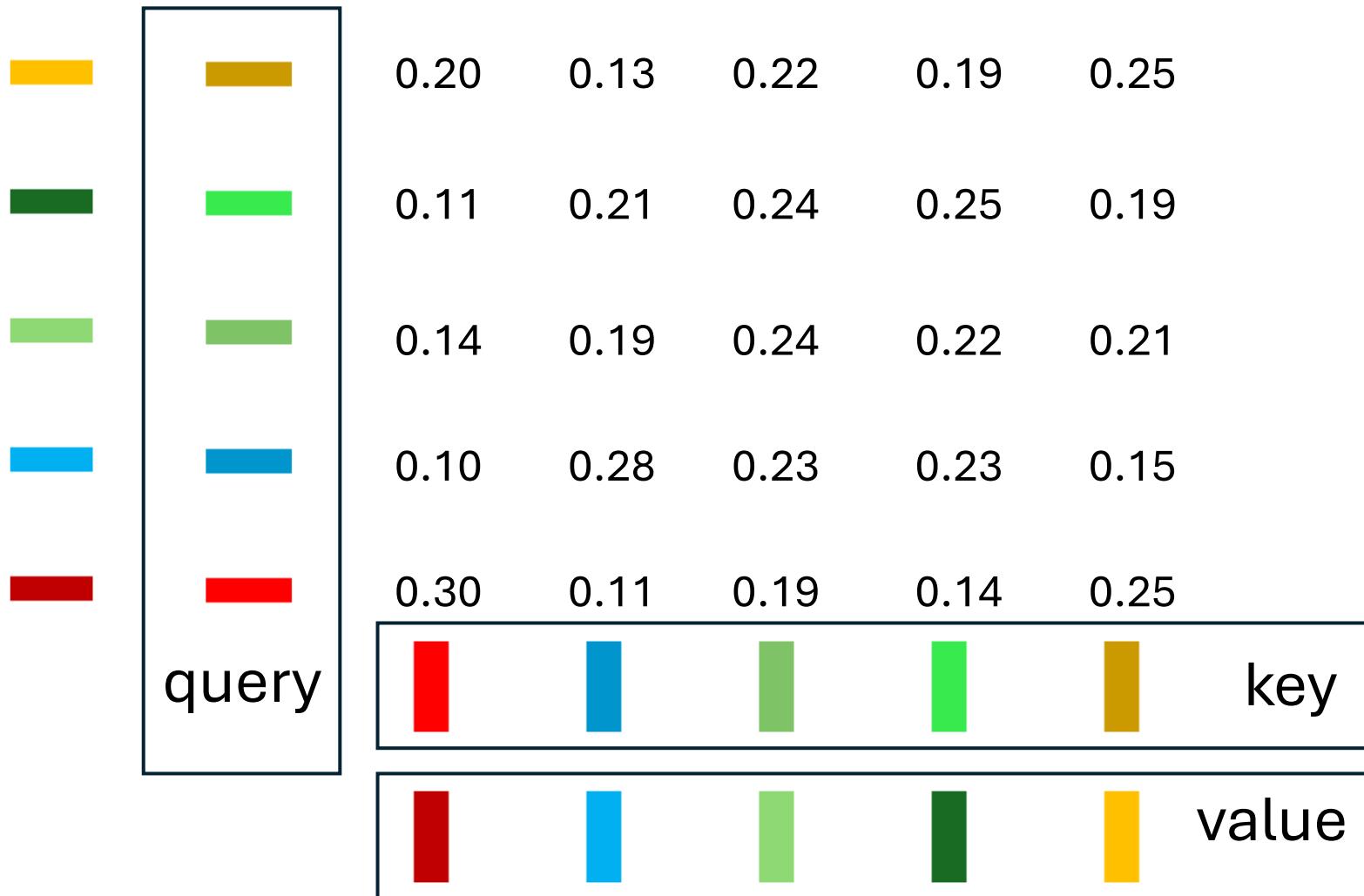
# Self-Attention Layer



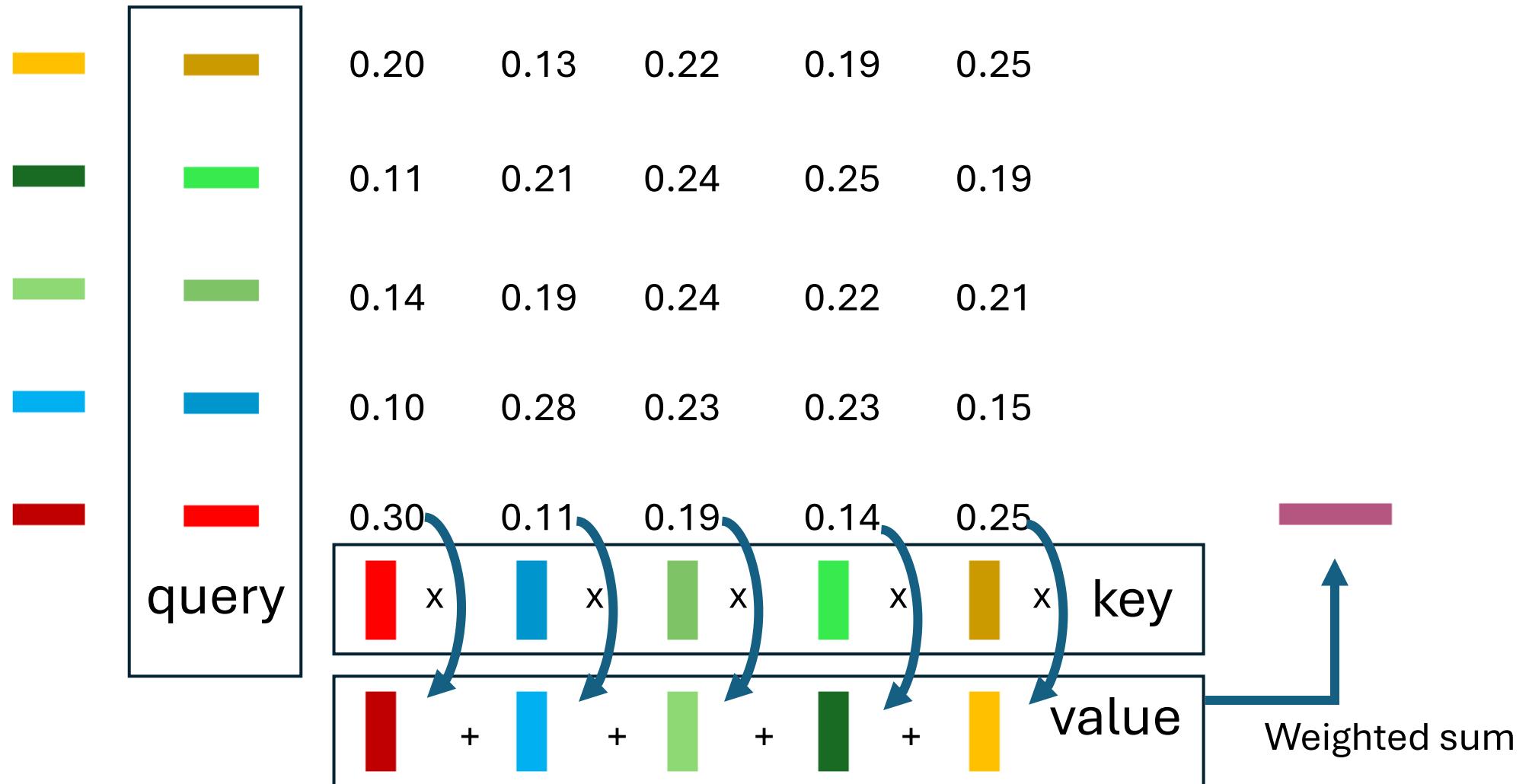
# Self-Attention Layer



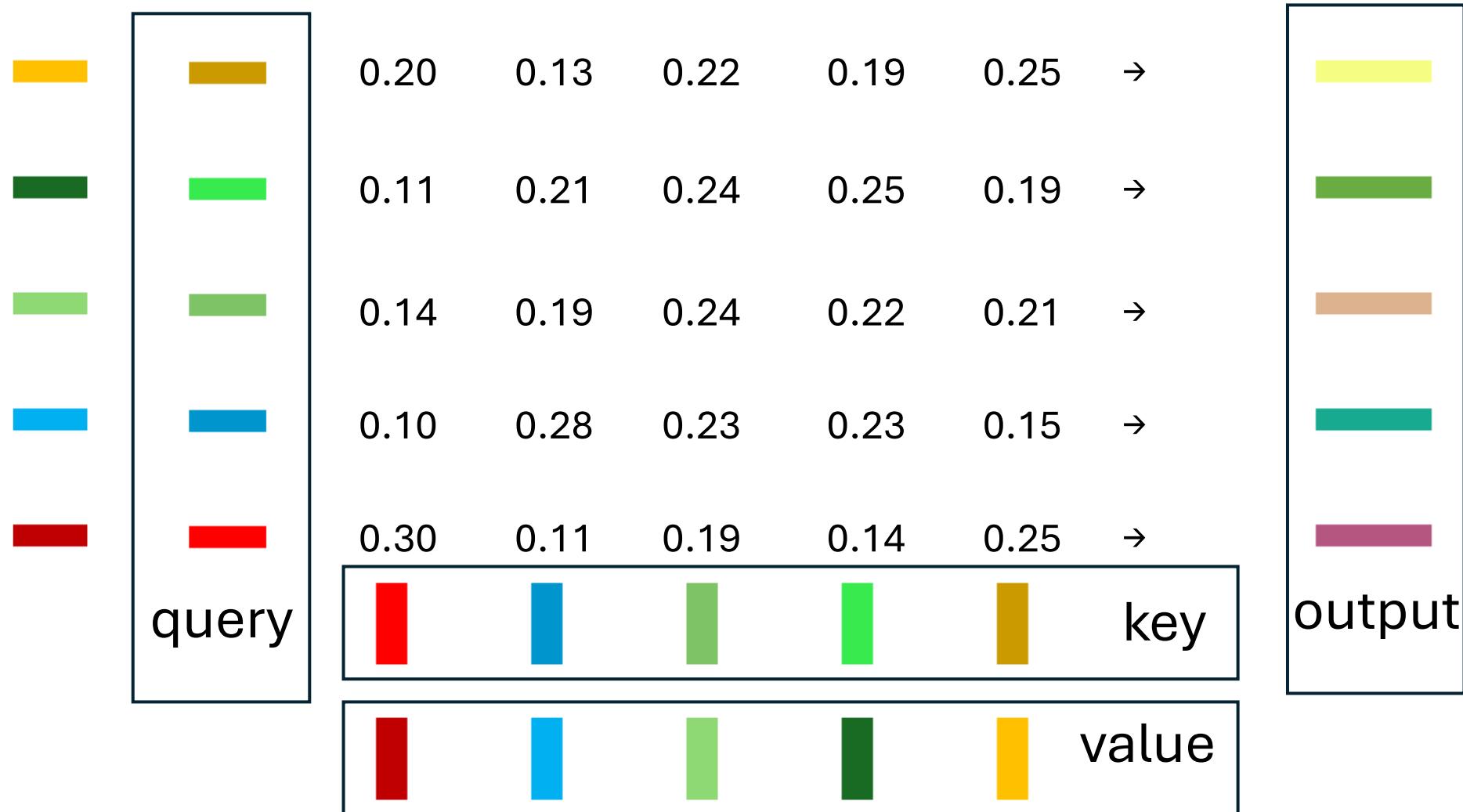
# Self-Attention Layer



# Self-Attention Layer



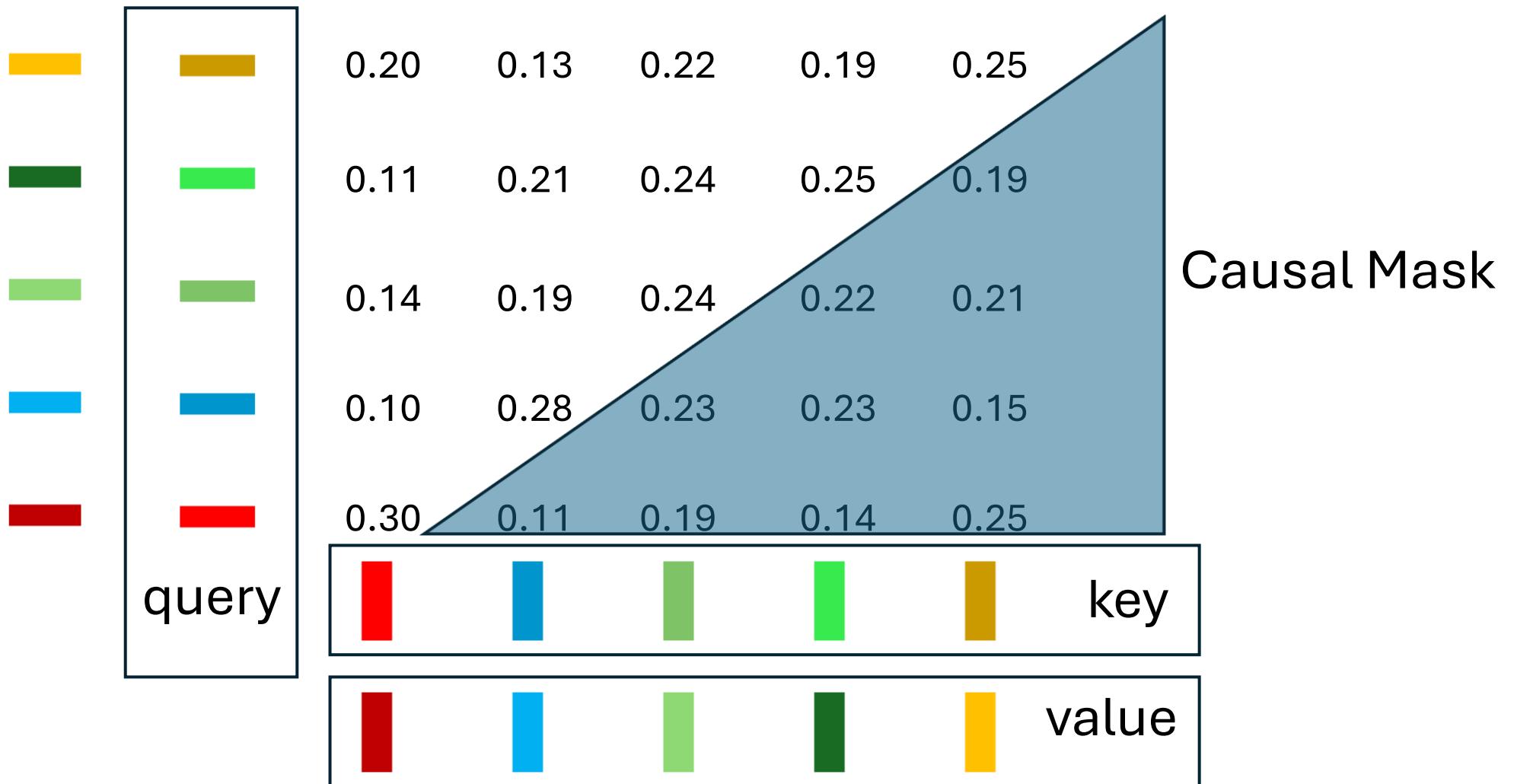
# Self-Attention Layer



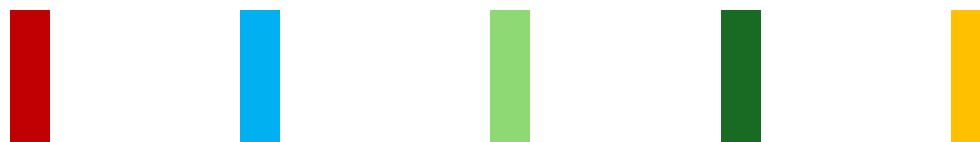
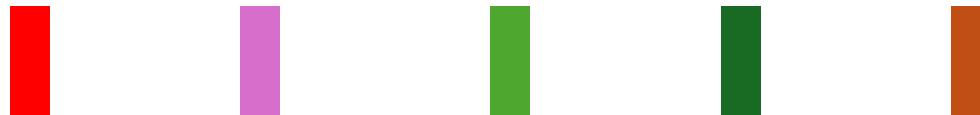
# Variations

- Attention with **Causal Mask**
  - Force left->right generation (also called autoregressive)
  - Useful for text generation)
- **Cross-attention** (useful to combine information from different modalities)
- **Multi-Head Attention (MHA)**
  - Use sets of Attention in parallel
  - Each Attention head focus on a particular features similar to Convolutional filter
  - **This is interesting. LSTM has units, Transformer has attention heads**
- Attention block = transform a sequence of vectors into another sequence of vectors of the same length
- **Positional Encoding** = Embedding layer for position of each token in input sequence

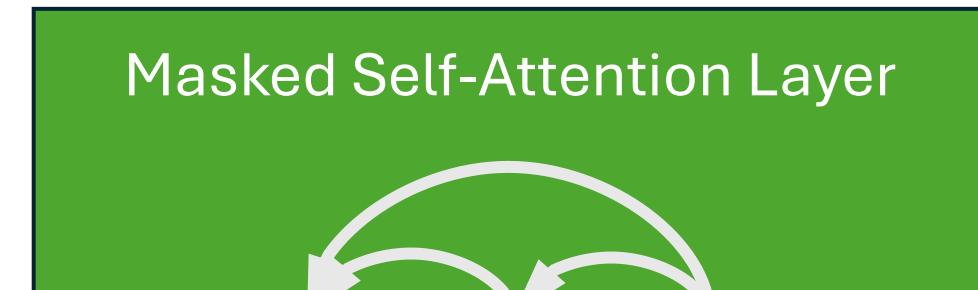
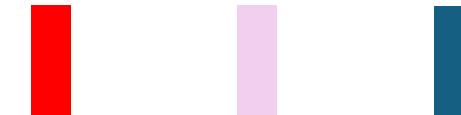
# Self-Attention Layer



## Transformer Encoder BERT



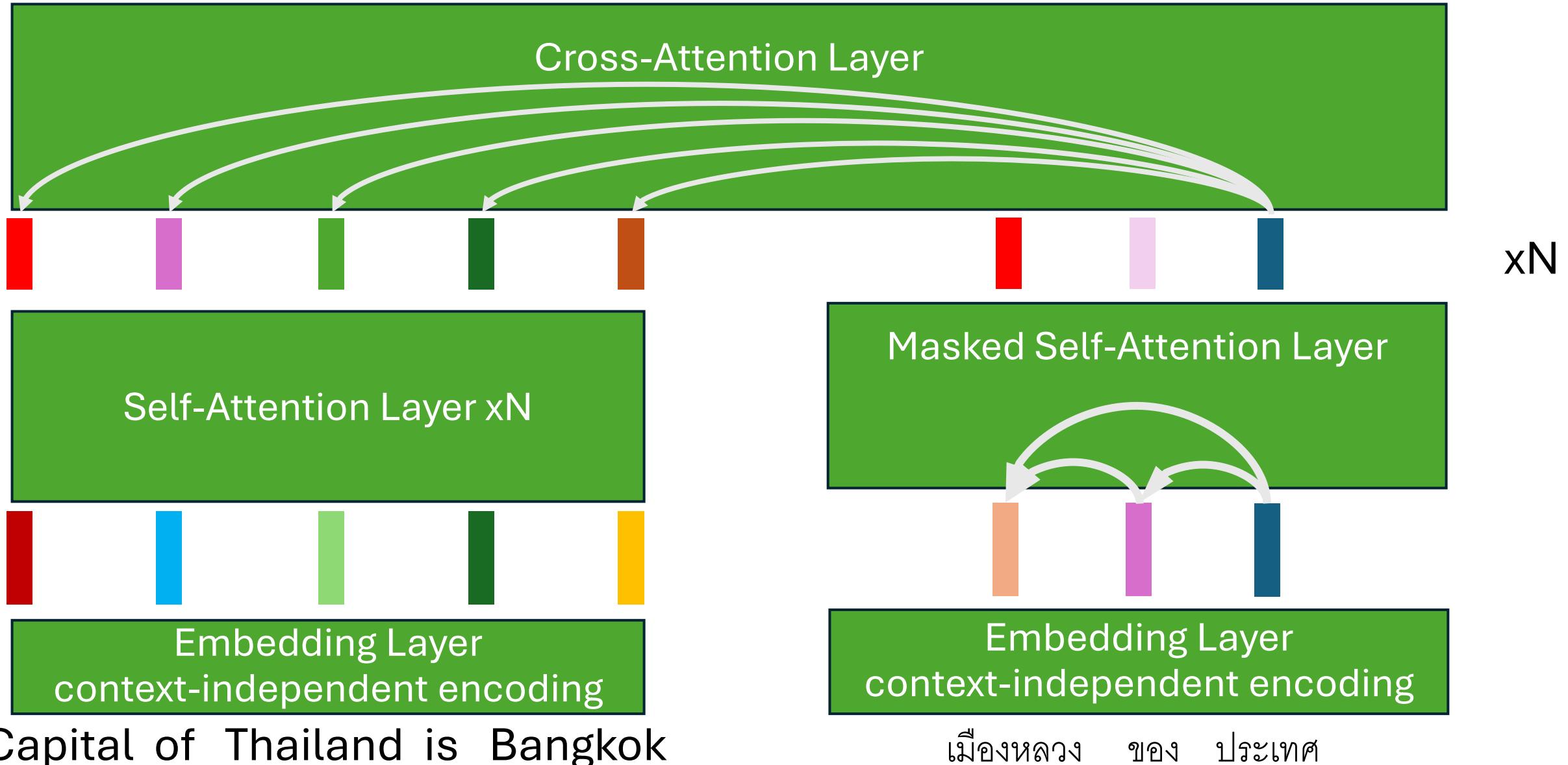
## Transformer Decoder

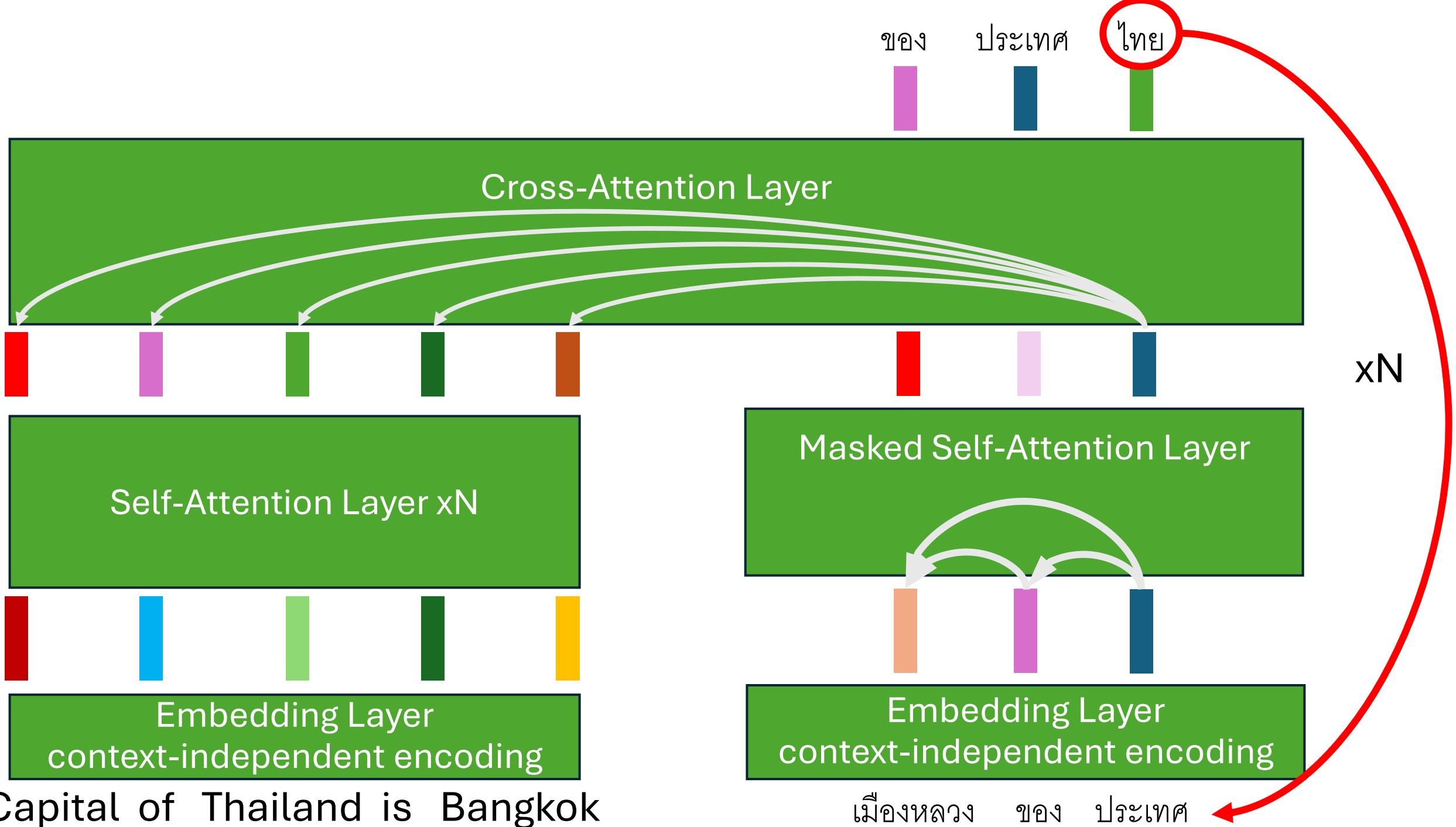


เมืองหลวง ของ ประเทศไทย

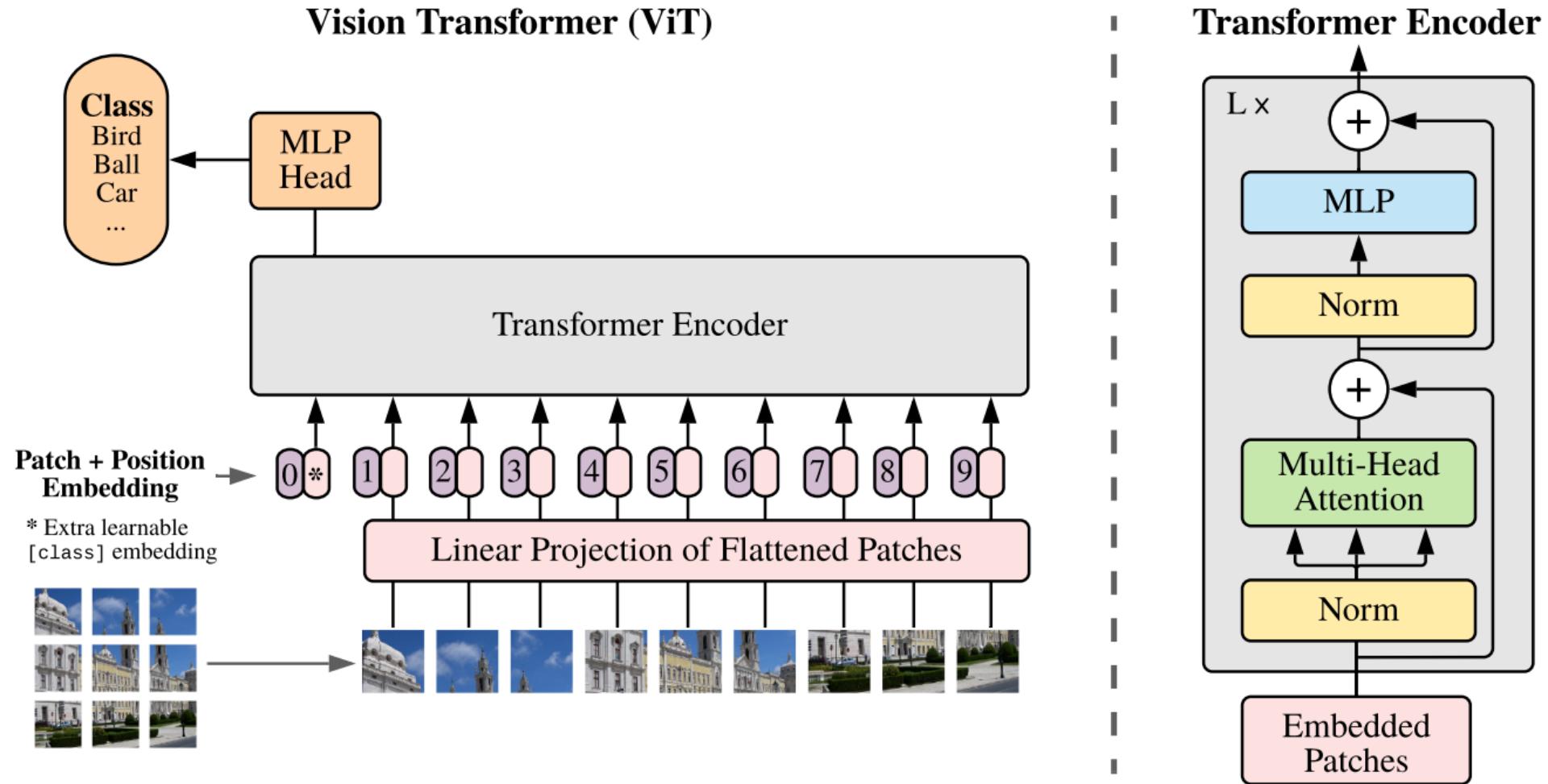
Capital of Thailand is Bangkok

ของ ประเทศ ไทย





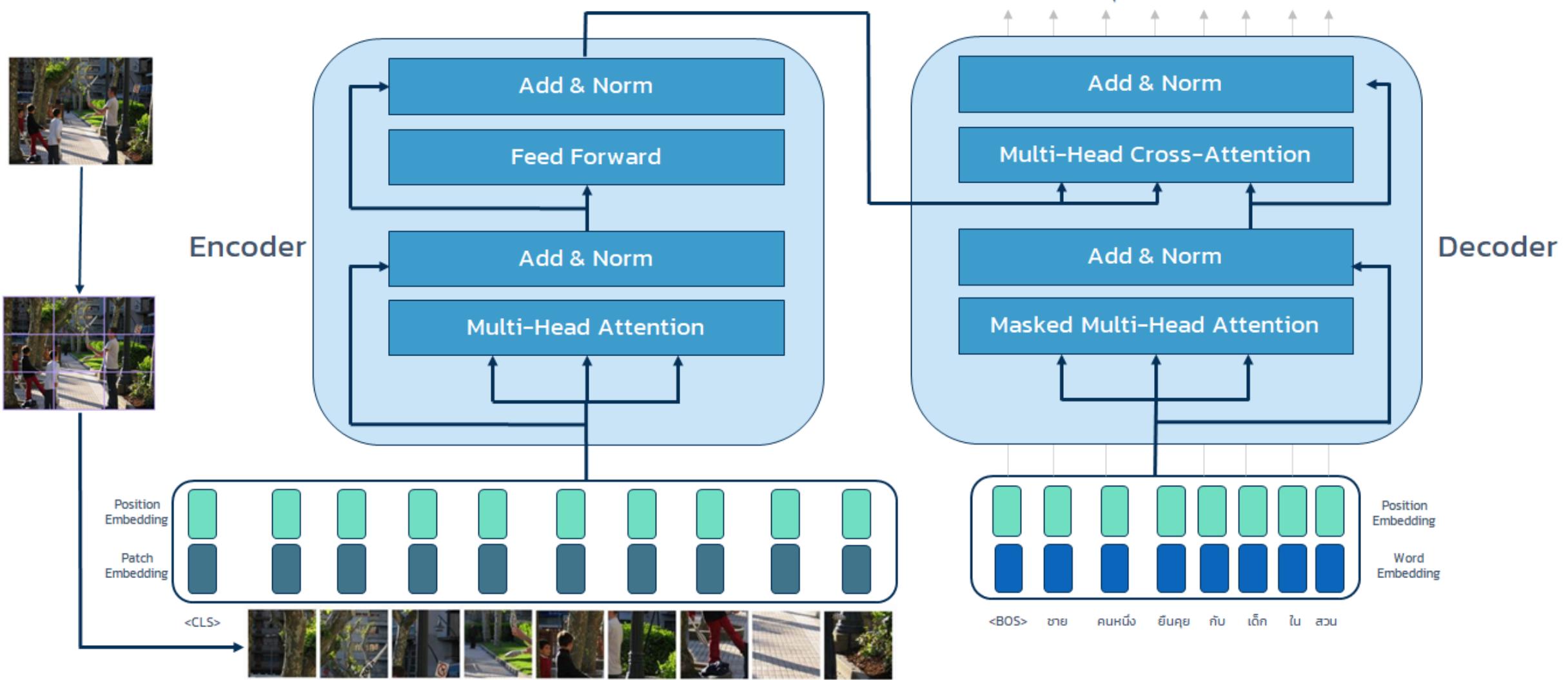
# Vision Transformer (ViT)

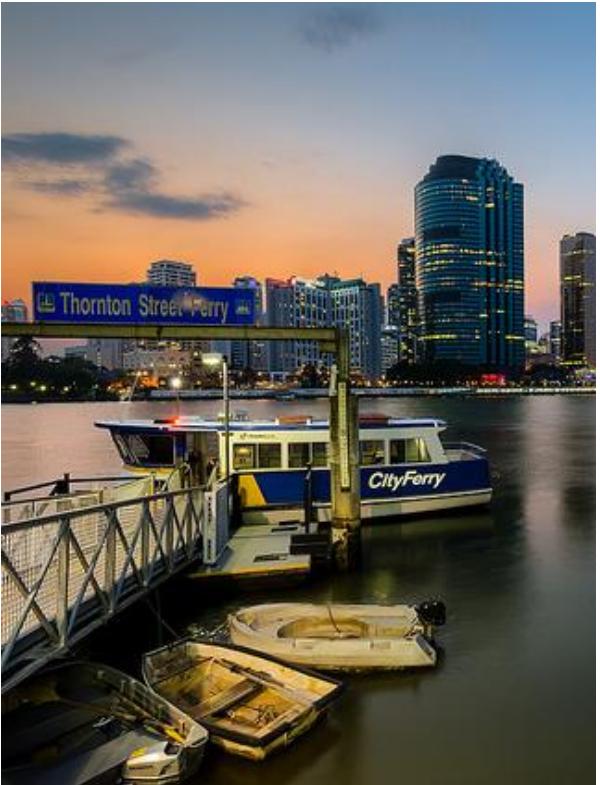


# Alternative approaches

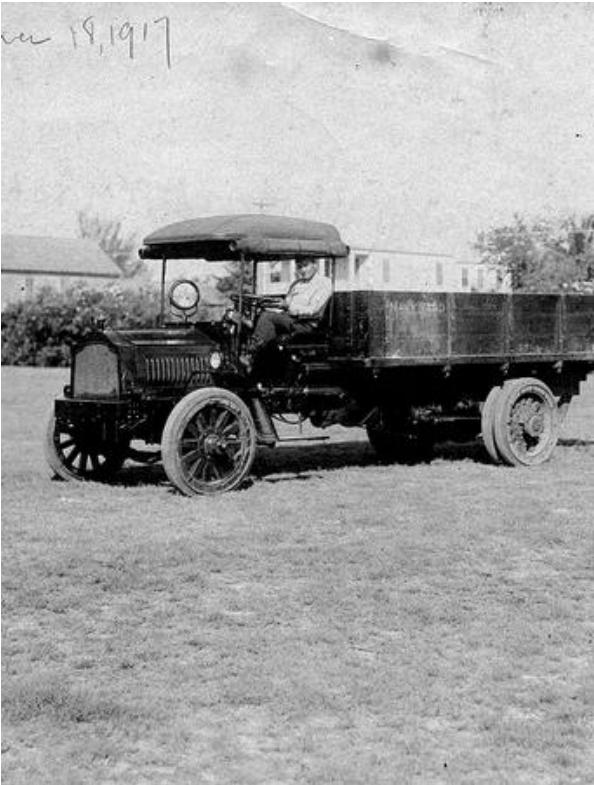
- Depth-wise Separable Convolution
- ConvNeXt
- CoAtNet
- MLP-Mixer
- Graph NN

# Image Captioning





แม่น้ำแห่งหนึ่งมีเรือจำนวน  
หลายลำจอดอยู่บันแม่น้ำ  
ด้านข้างสะพาน



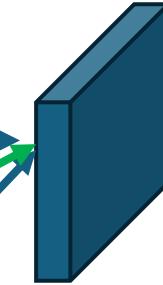
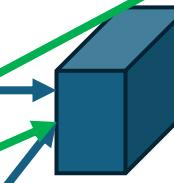
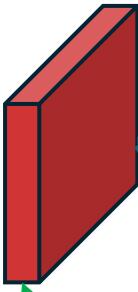
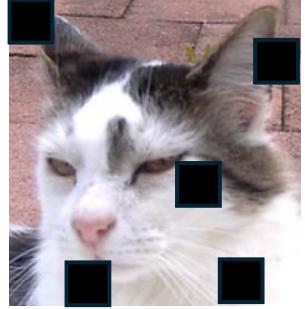
สนามหญ้ามีรถบรรทุกคัน  
สีดำคันหนึ่งจอดอยู่ข้างบ้าน  
หลังหนึ่ง



คนกำลังเล่นว่าวอยู่ใกล้กับ  
ตึกสูงสีขาวและห้องพักมี  
เมฆมาก



ผู้ชายผอมสีดำมีหนวดสวม  
แว่นกำลังถือลูกบอลงสี  
เหลือง

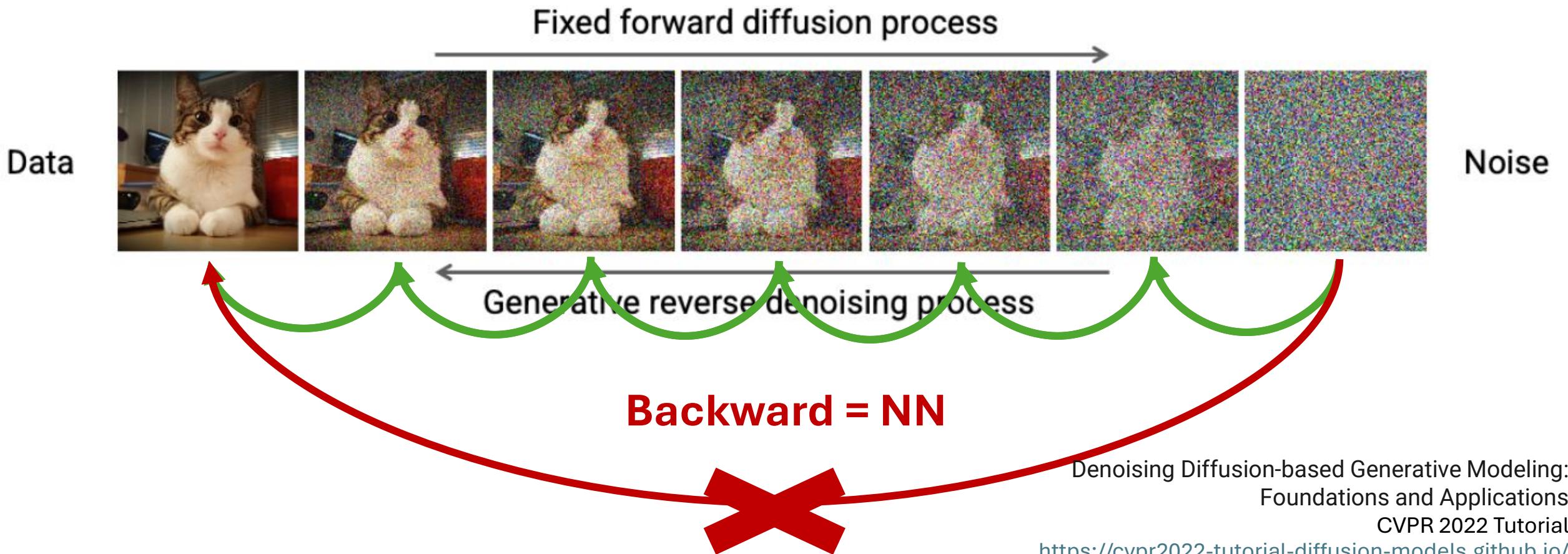


An image of a cat

Stable Diffusion

# Diffusion Model

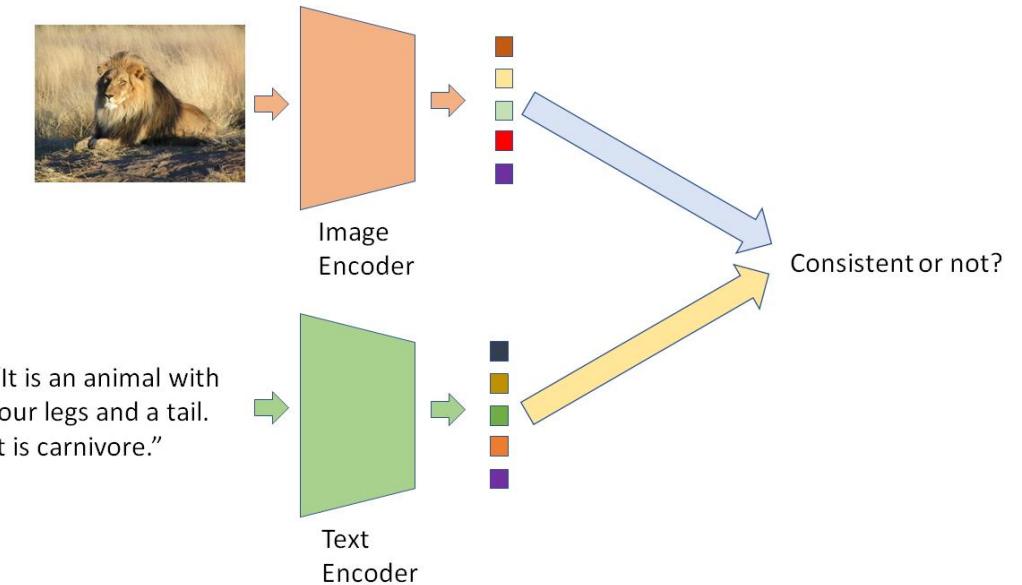
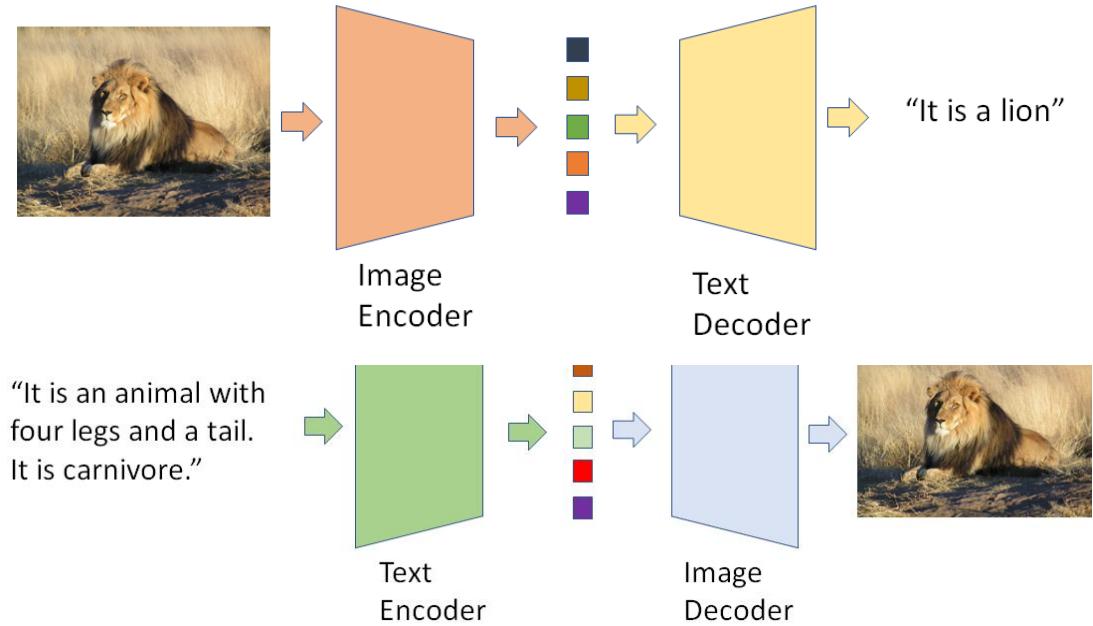
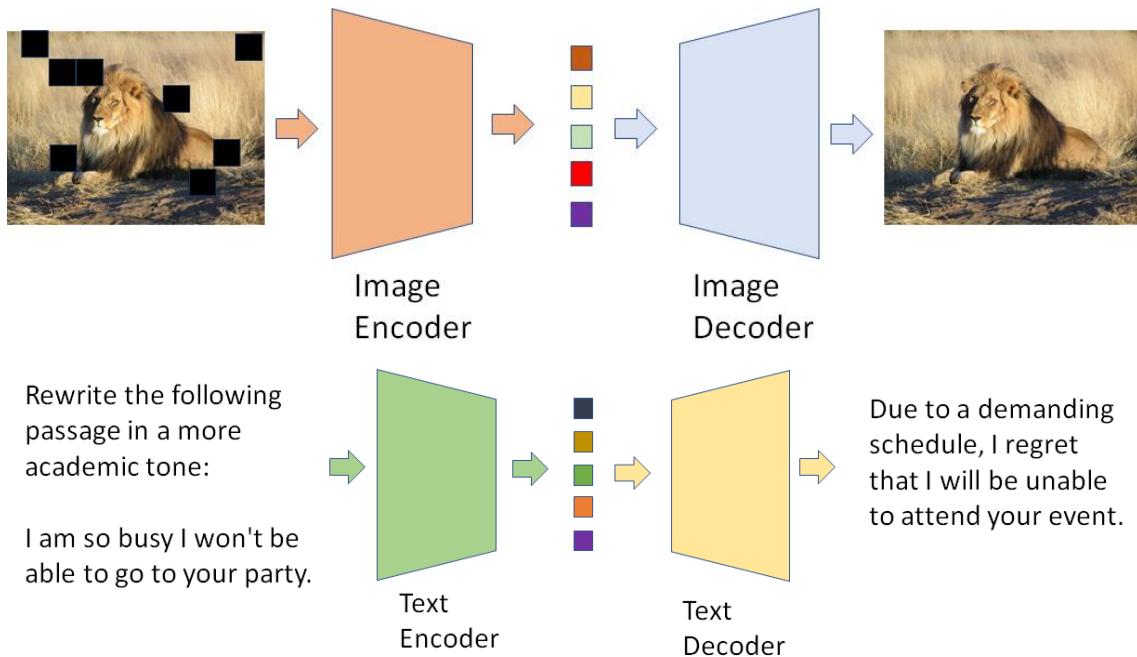
**Forward = Add Noise**





# Future of Computer Vision

# Multi-Modal AI



# Image Search

- Use surrounding text as index of an image in a web page work with text query
- Use image feature vector as index of an image, work with image query (Content-Based Image Retrieval, CBIR)
- Use object detectors, e.g. YOLO, to index images with known keywords work with text query
- Use CLIP visual feature vector as index of an image
  - Compare against CLIP text feature vector from text query

- Synthetic Data
- Federated Learning
- Out-of-distribution
- Frameworks: Tensorflow Lite, TinyML
- Model:
  - Quantization, Pruning
  - Neural Architecture Search (NAS)
- Spiking NN
- Explanable AI (XAI)
- Counterfactual Analysis

# Thank you

Q & A