

# Lecture 10:

# Open Domain Question Answering, Speech recognition and Distillation

MADE, Moscow

03.06.2020

**Radoslav Neychev**

1. Open domain Question Answering
  - DrQA
2. Speech recognition and generation: brief overview
  - ASR
3. Extra: Knowledge Distillation
4. NLP area and perspectives overview
5. Outro

Based on: [Series of Medium posts](#)

<http://ai.stanford.edu/blog/answering-complex-questions/>  
[Distilling the Knowledge in a Neural Network](#)

# Open Domain question answering

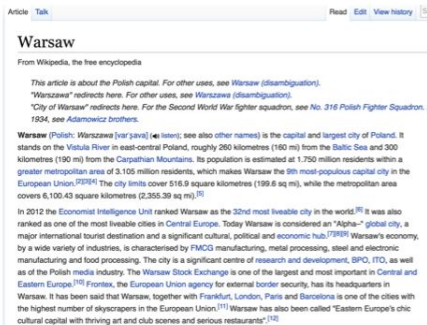
## Open-domain QA

SQuAD, TREC, WebQuestions, WikiMovies

Q: How many of Warsaw's inhabitants spoke Polish in 1933?

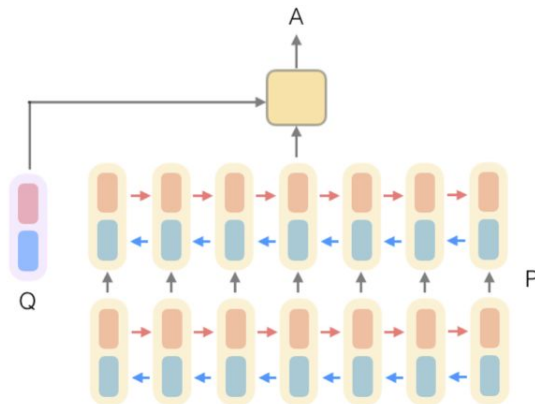


## Document Retriever



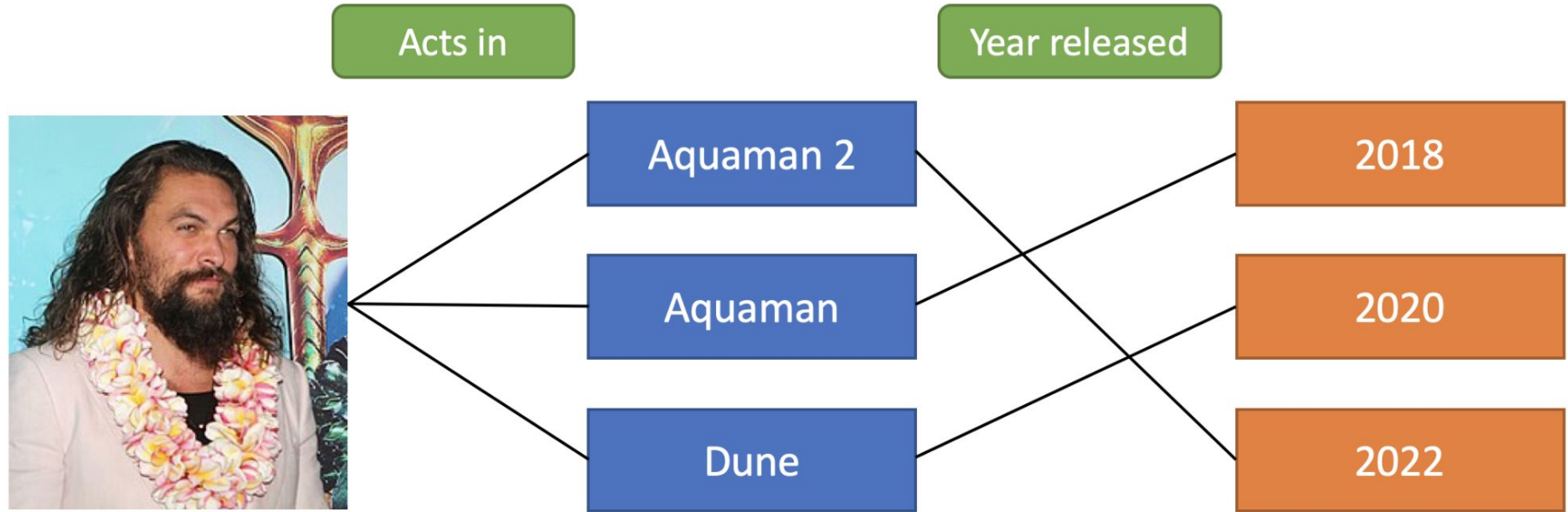
# Document Reader

833,500

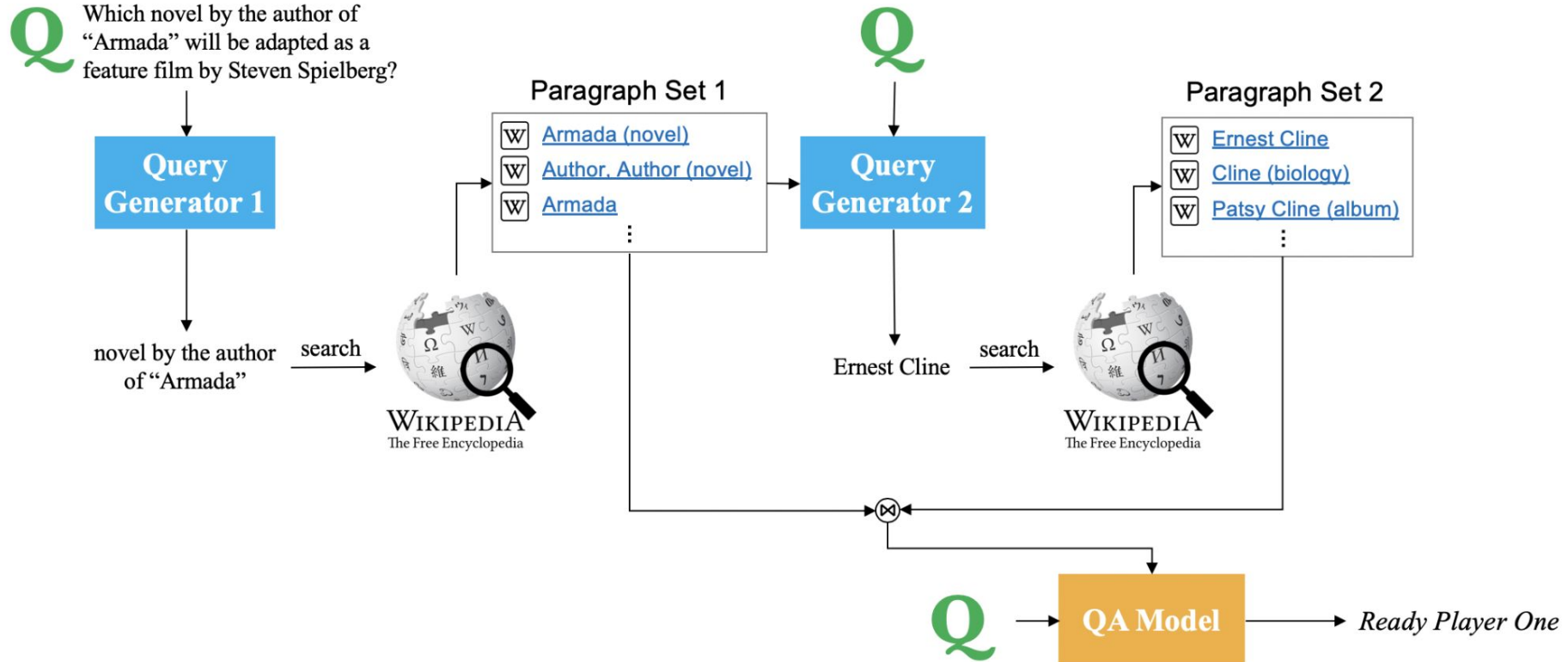


# Possible problems

*Example question: “What is the Aquaman actor’s next movie?”*



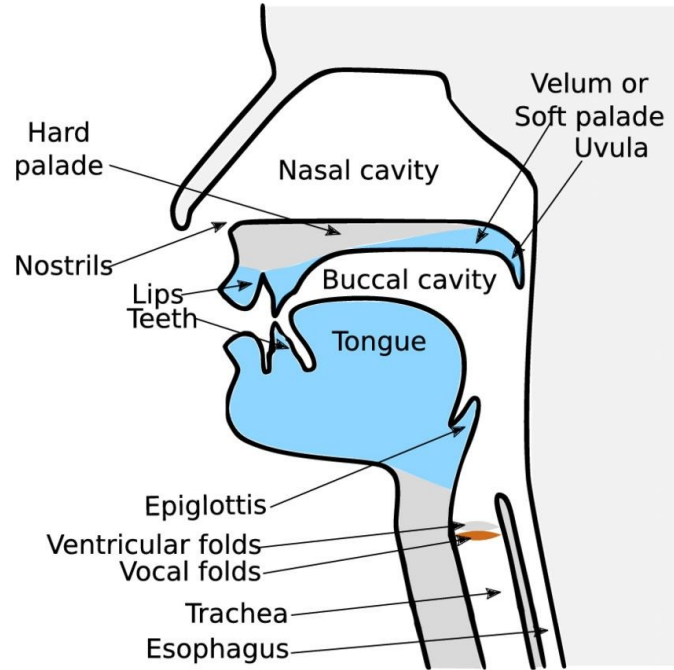
# Possible solutions



# Speech recognition and generation

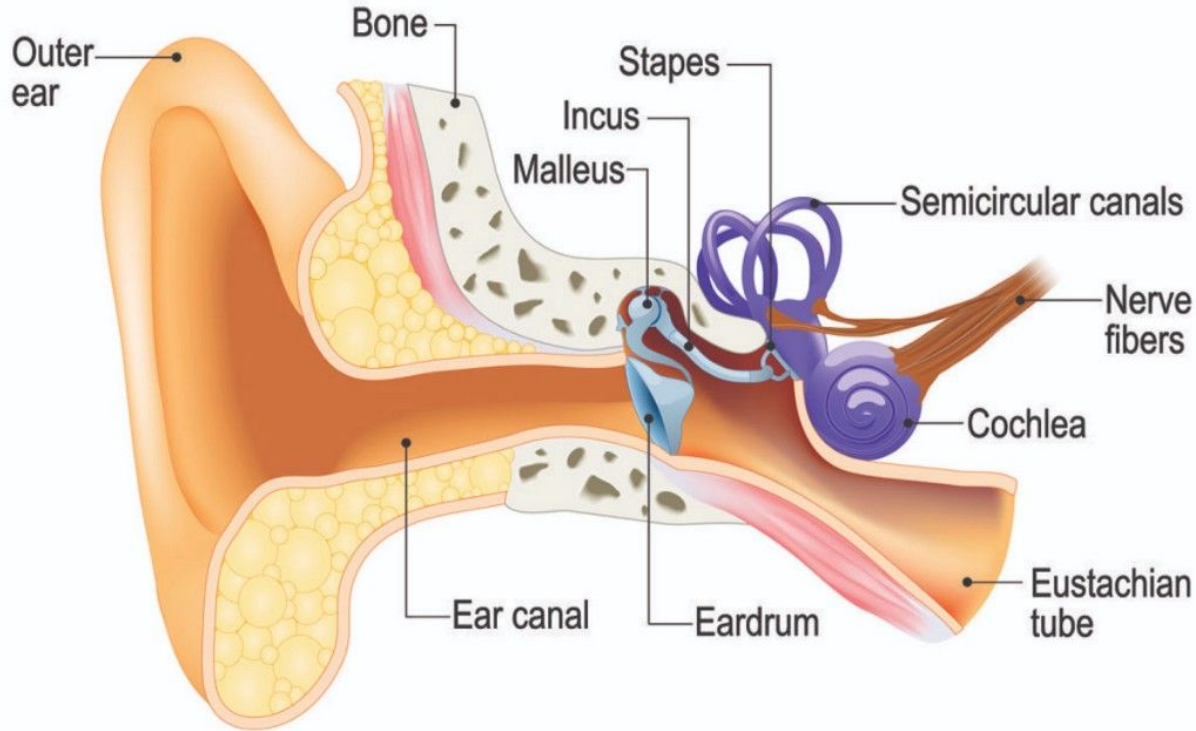
## ASR and TTS

# Foundations: vocal tract





# Foundations: ear structure



# Foundations: ear structure

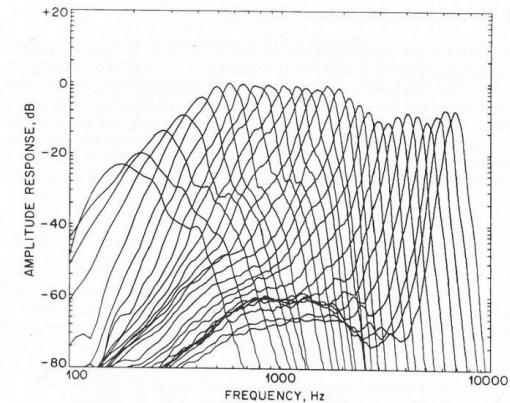
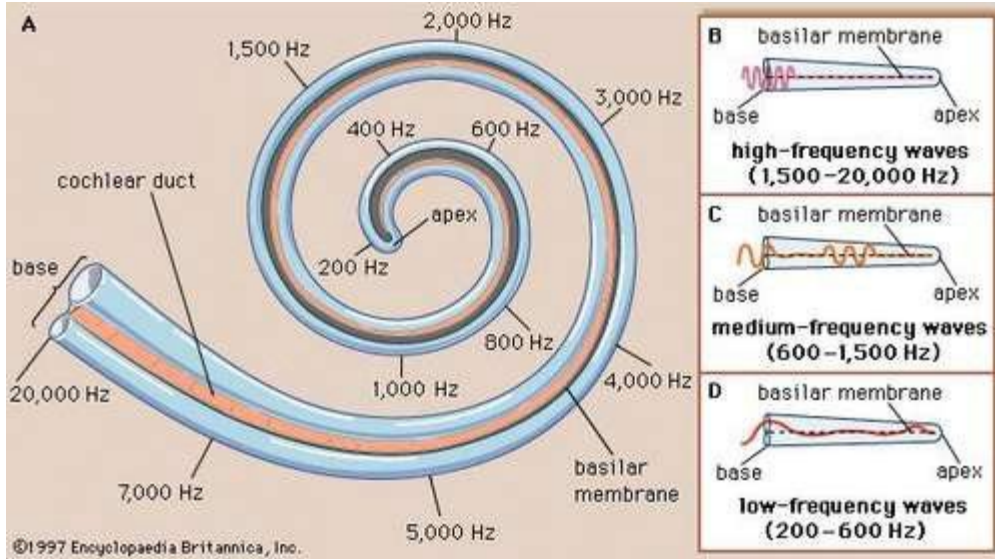


Figure 3.50 Frequency response curves of a cat's basilar membrane (after Ghitza [13]).

# Processing the sound

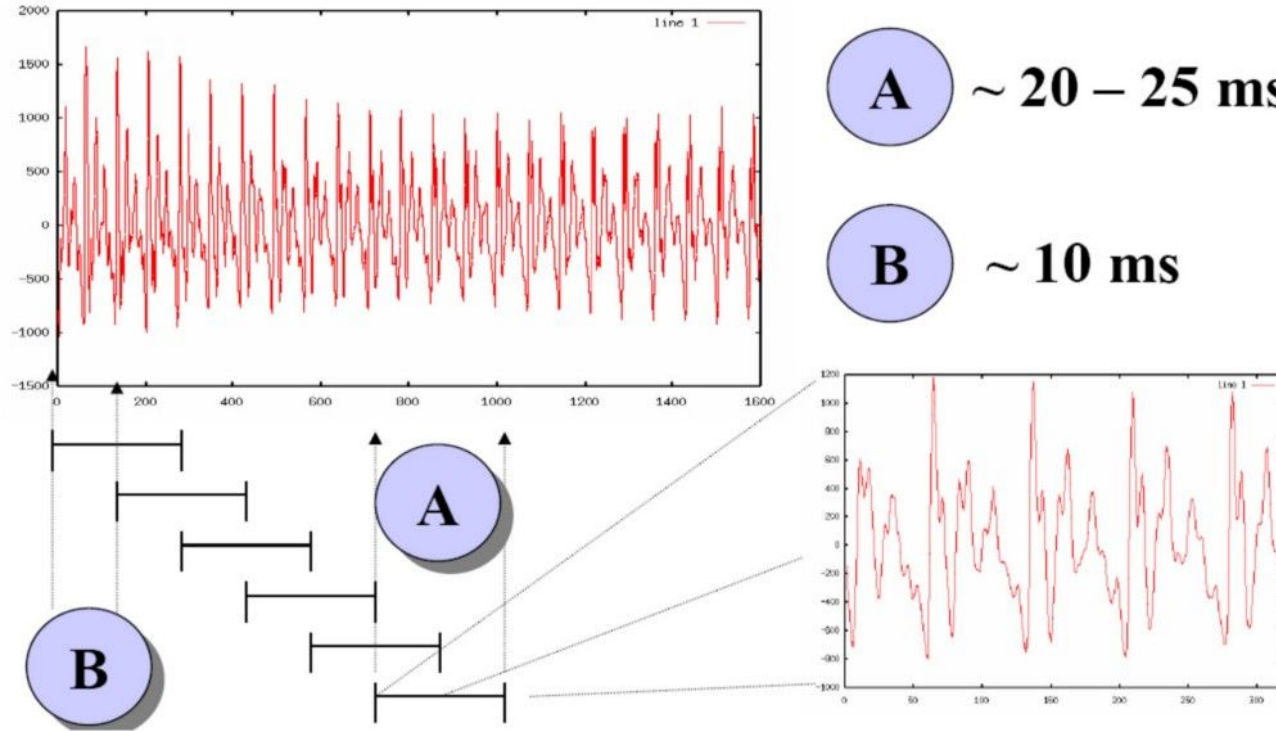


Image from Bryan Pellom

# Fourier Transform

## Definition [\[ edit \]](#)

The *discrete Fourier transform* transforms a [sequence](#) of  $N$  [complex numbers](#)

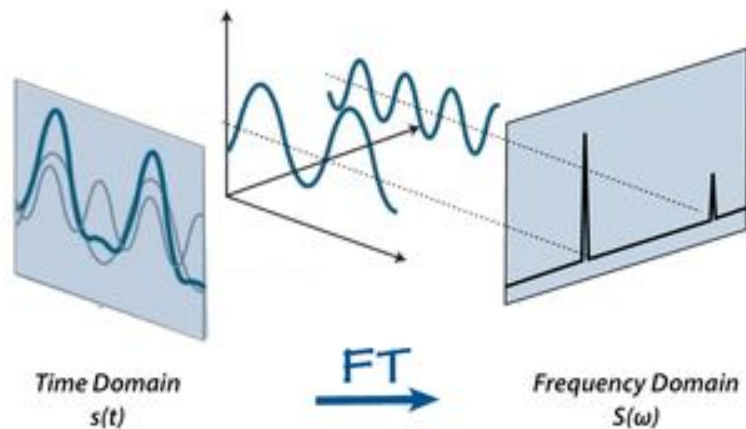
$\{\mathbf{x}_n\} := x_0, x_1, \dots, x_{N-1}$  into another sequence of complex numbers,

$\{\mathbf{X}_k\} := X_0, X_1, \dots, X_{N-1}$ , which is defined by

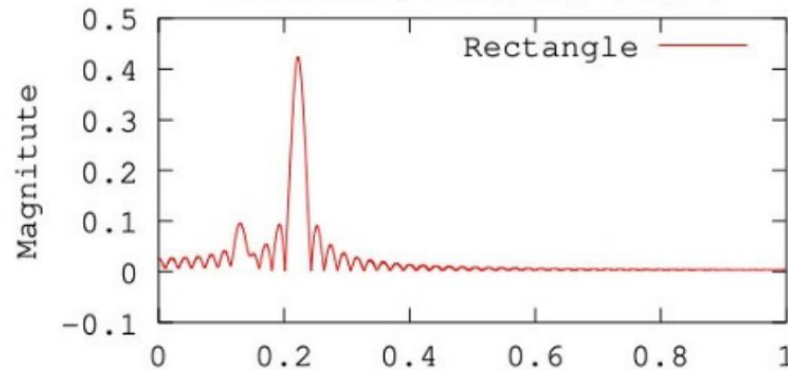
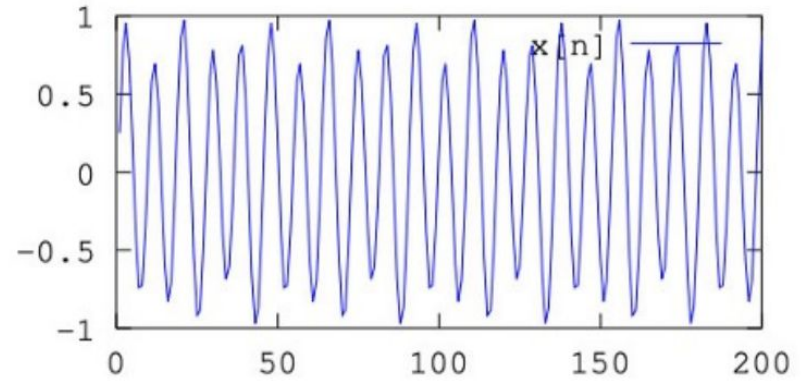
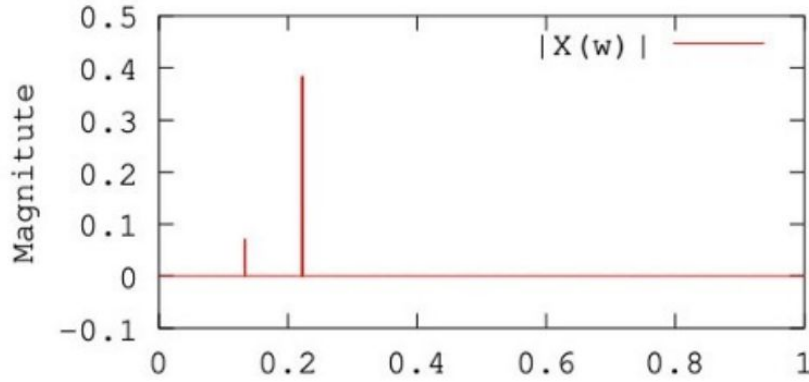
$$\begin{aligned} X_k &= \sum_{n=0}^{N-1} x_n \cdot e^{-\frac{i2\pi}{N}kn} \\ &= \sum_{n=0}^{N-1} x_n \cdot \left[ \cos\left(\frac{2\pi}{N}kn\right) - i \cdot \sin\left(\frac{2\pi}{N}kn\right) \right], \end{aligned} \quad (\text{Eq.1})$$

where the last expression follows from the first one by [Euler's formula](#).

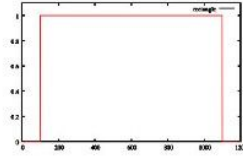
The transform is sometimes denoted by the symbol  $\mathcal{F}$ , as in  $\mathbf{X} = \mathcal{F}\{\mathbf{x}\}$  or  $\mathcal{F}(\mathbf{x})$  or  $\mathcal{F}\mathbf{x}$ .<sup>[A]</sup>



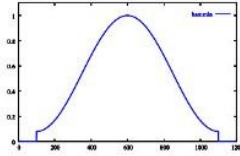
# Fourier Transform: filtering



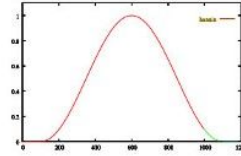
# Fourier Transform: filtering



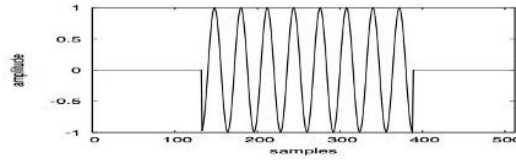
Rectangular



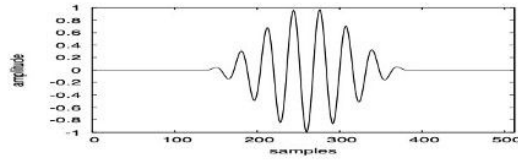
Hamming



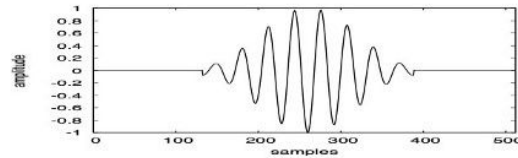
Hanning



(a) Rectangular window



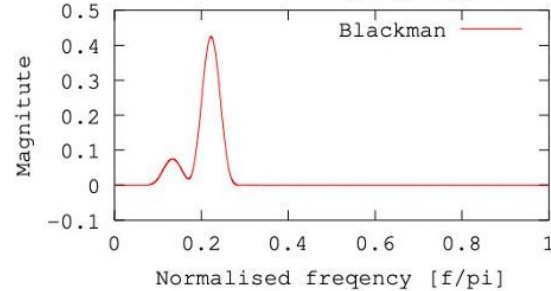
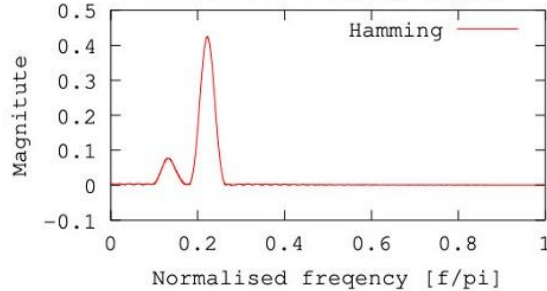
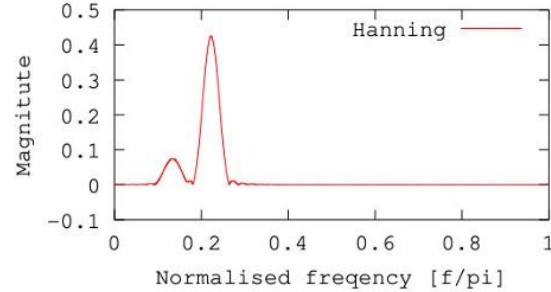
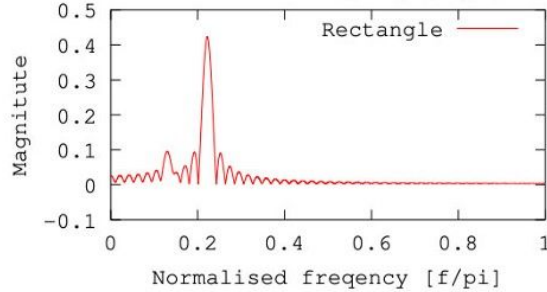
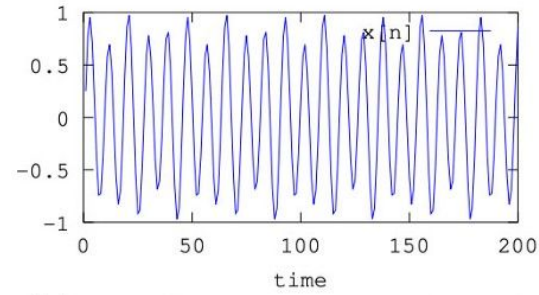
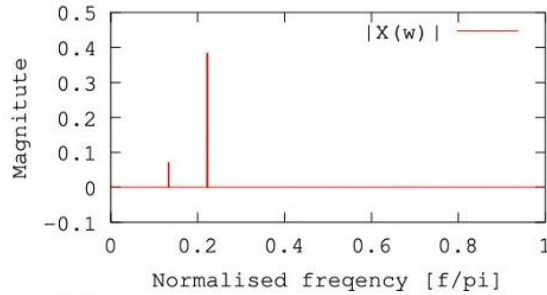
(b) Hanning window



(c) Hamming window

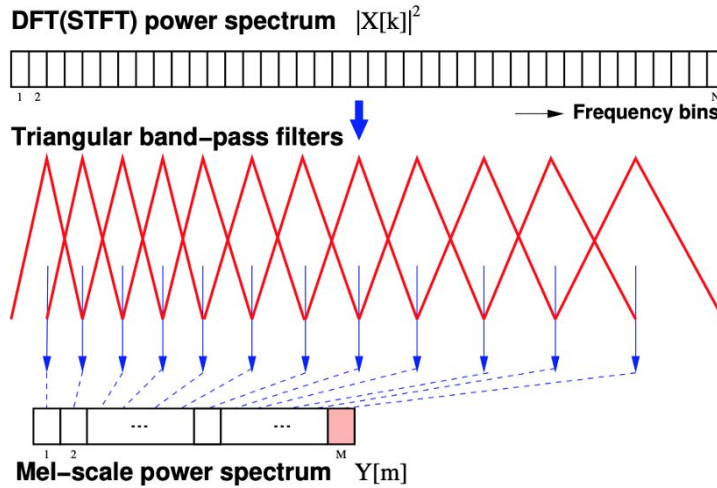
(Taylor, fig 12.1)

# Fourier Transform: filtering





# Features generation

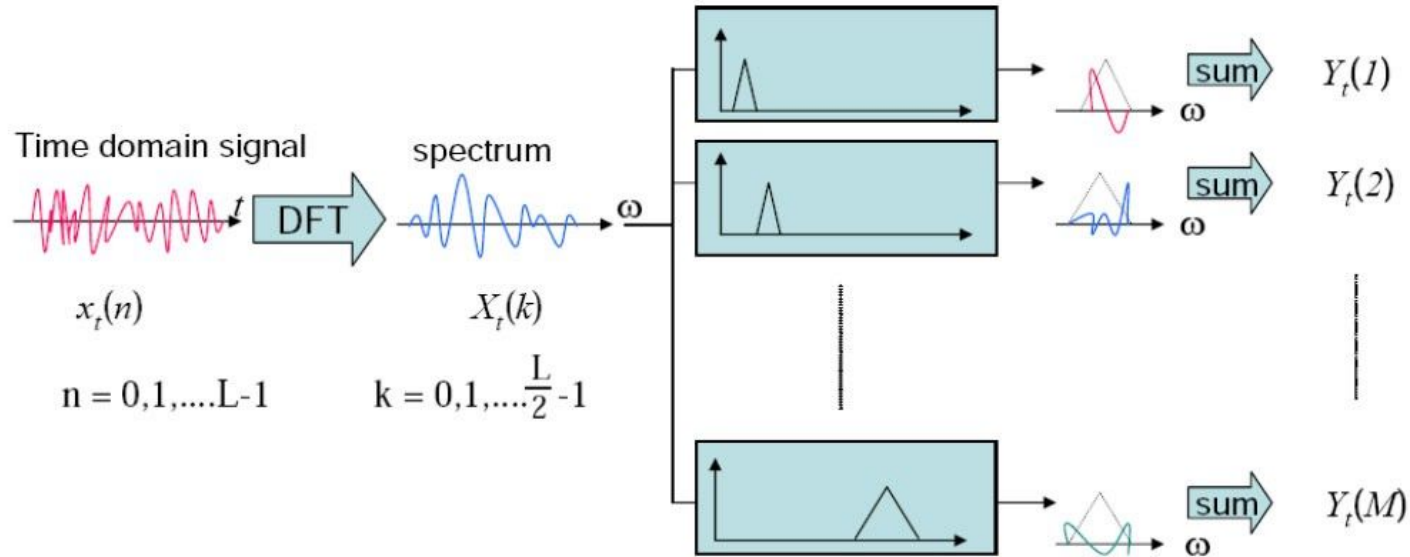


$$Y_t[m] = \sum_{k=1}^N W_m[k] |X_t[k]|^2$$

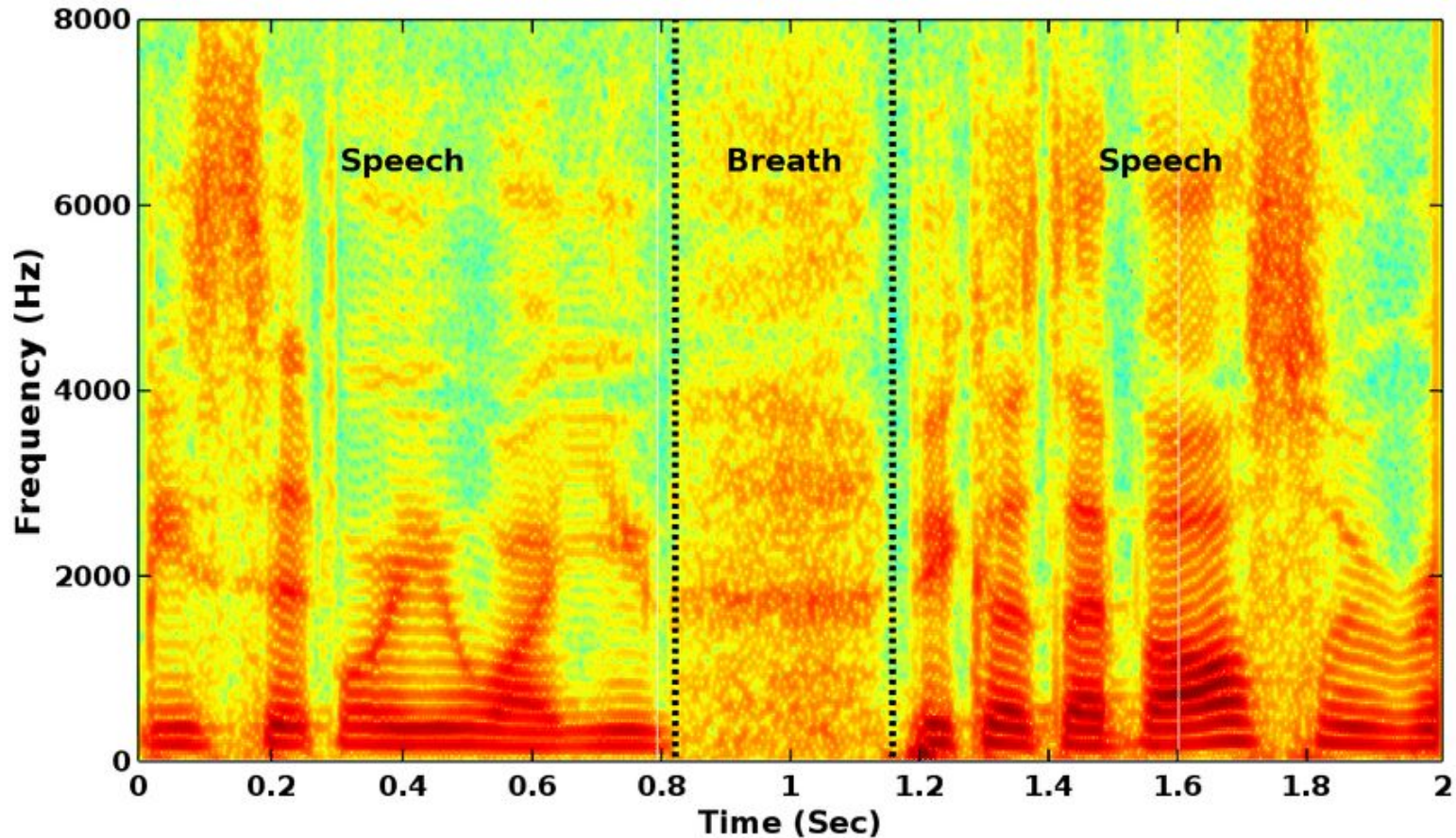
where  $k$  : DFT bin number  $(1, \dots, N)$   
 $m$  : mel-filter bank number  $(1, \dots, M)$



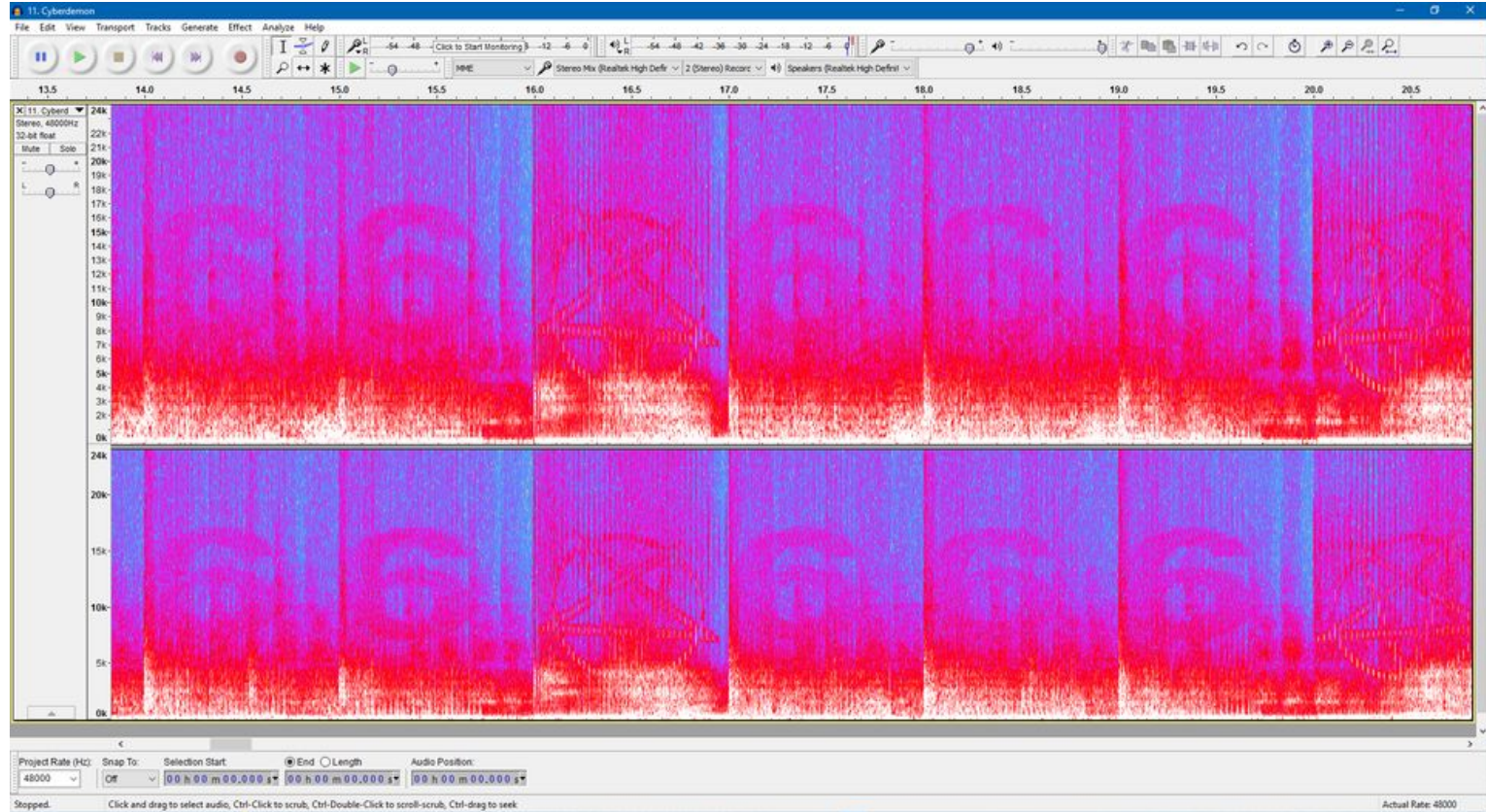
# Features generation



# Spectrogram



# Spectrograms as Easter eggs. DOOM 2016 soundtrack



# Embeddings

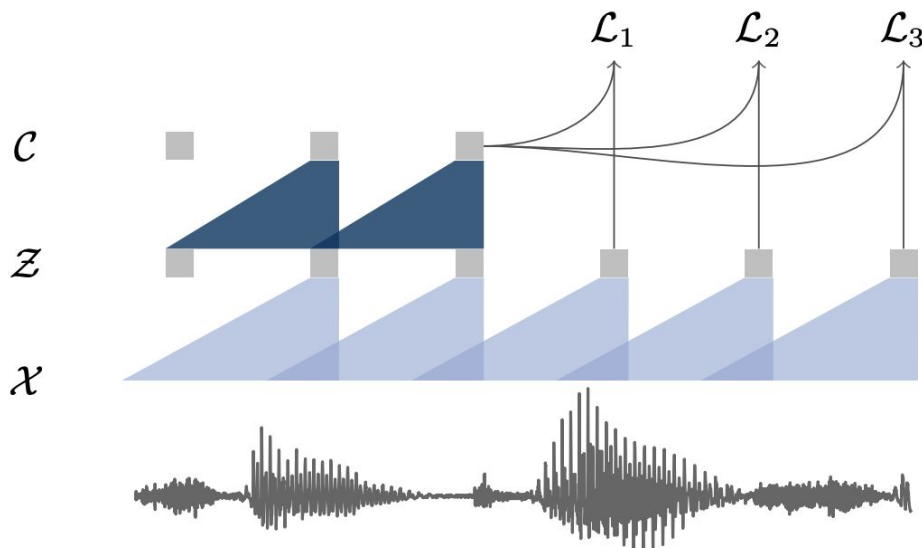


Figure 1: Illustration of pre-training from audio data  $\mathcal{X}$  which is encoded with two convolutional neural networks that are stacked on top of each other. The model is optimized to solve a next time step prediction task.

# Extra: Knowledge Distillation

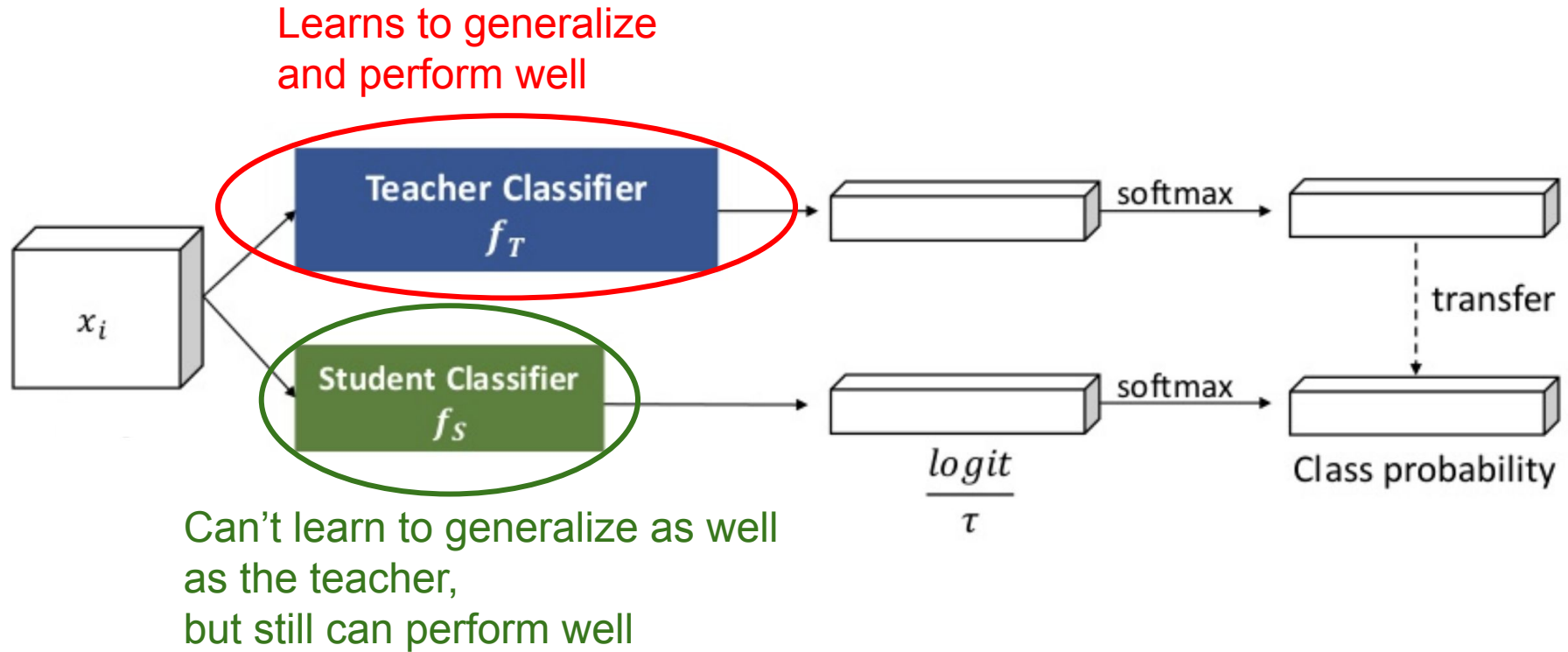


# Cerura Vinula in caterpillar and butterfly forms



Do they have the same “life purpose”  
and solve the same problems?

# Knowledge distillation



# Knowledge distillation

Denote **teacher** and **student** models.

**Student** model has logits  $z_i$  and corresponding probabilities  $q_i$ , derived with the softmax operation:

$$q_i = \frac{\exp(z_i/T)}{\sum_j \exp(z_j/T)}$$

where  $T$  stays for the temperature.

**Teacher** model has logits  $v_i$  and corresponding probabilities  $p_i$ .



# Knowledge distillation

Let's derive the cross-entropy gradient on **student** logits using the **teacher** predictions as targets:

$$\frac{\partial C}{\partial z_i} = \frac{1}{T} (q_i - p_i) = \frac{1}{T} \left( \frac{e^{z_i/T}}{\sum_j e^{z_j/T}} - \frac{e^{v_i/T}}{\sum_j e^{v_j/T}} \right)$$

If the temperature is high, the following equation takes place:

$$\frac{\partial C}{\partial z_i} \approx \frac{1}{T} \left( \frac{1 + z_i/T}{N + \sum_j z_j/T} - \frac{1 + v_i/T}{N + \sum_j v_j/T} \right)$$


# Knowledge distillation

Logits can be centered, so

$$\sum_j z_j = \sum_j v_j = 0$$

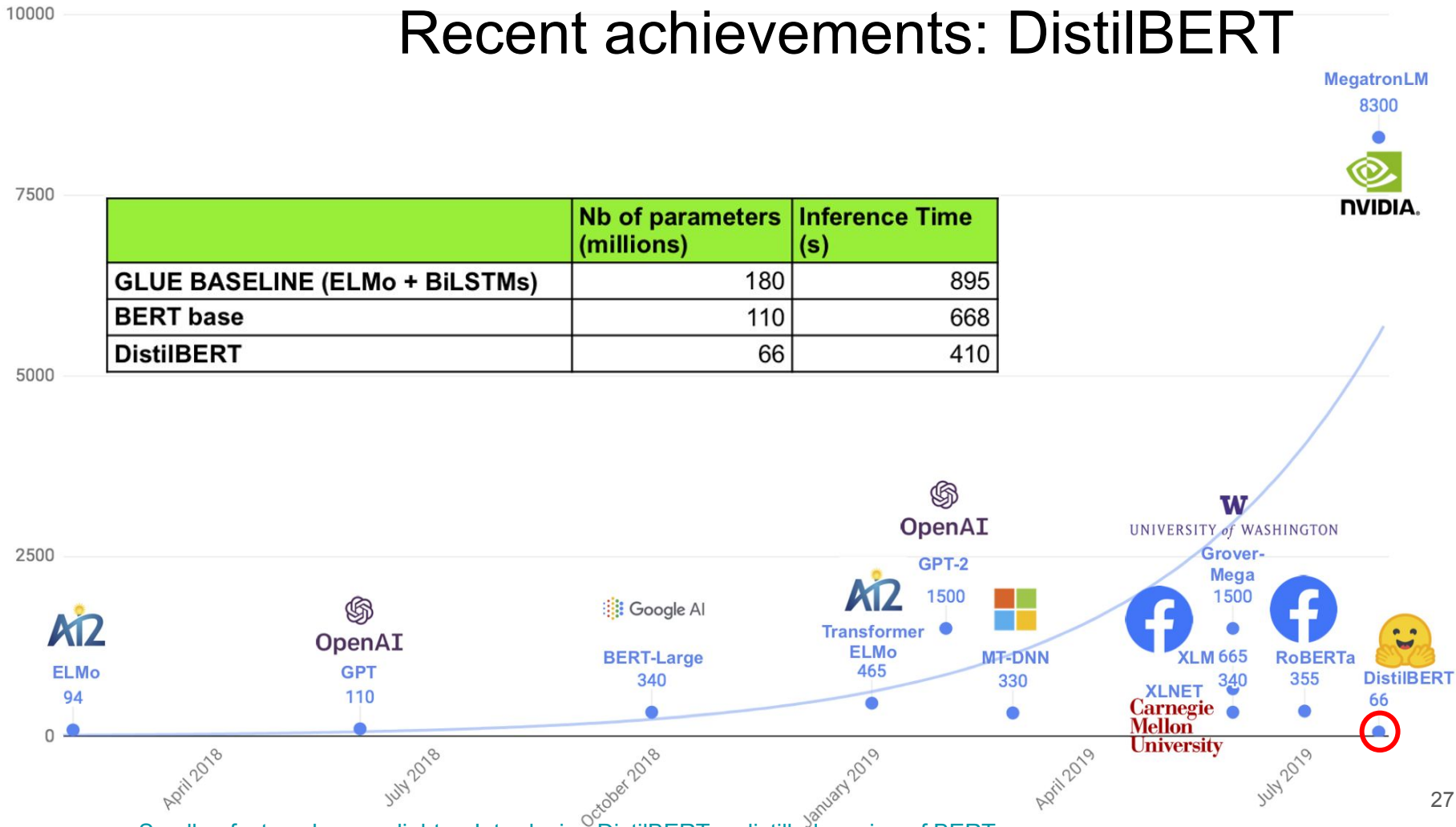
Then the gradient takes form:

$$\frac{\partial C}{\partial z_i} \approx \frac{1}{T} \left( \frac{1 + z_i/T}{N + \sum_j z_j/T} - \frac{1 + v_i/T}{N + \sum_j v_j/T} \right) \approx \frac{1}{NT^2} (z_i - v_i)$$

$$\frac{dC}{dz_i} = \frac{1}{NT^2} (z_i - v_i) \sim (z_i - v_i) = \overset{\text{Constant}}{M} \frac{d(z_i - v_i)^2}{dz_i}$$


# Recent achievements: DistilBERT

number of parameters, millions



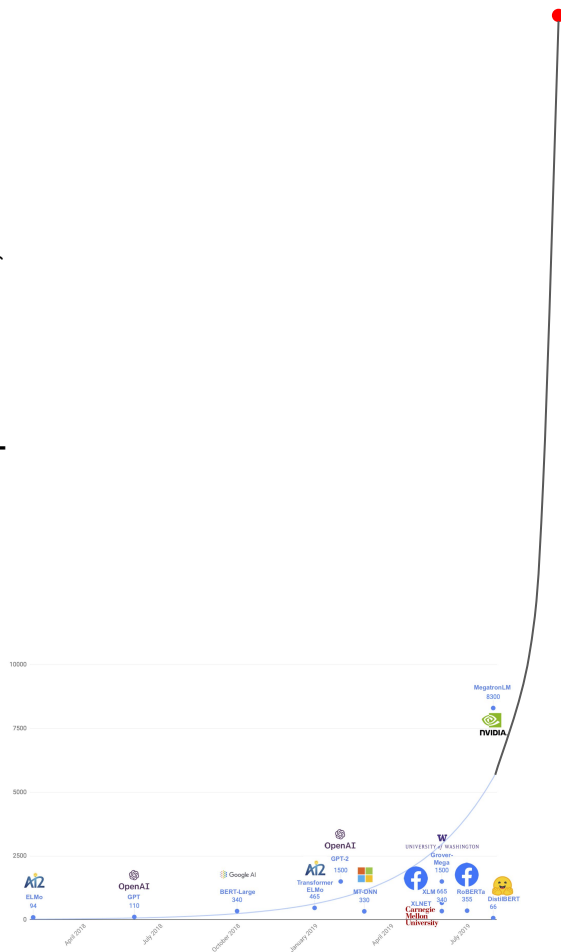
# Main ideas

- DistilBERT is initialized from its teacher, BERT, by taking one layer out of two, leveraging the common hidden size.
  - *Comment: Training a sub-network is not only about the architecture. It is also about finding the right initialization for the sub-network to converge.*
- DistilBERT is trained on very large batches leveraging gradient accumulation (up to 4000 examples per batch), with dynamic masking and removed the next sentence prediction objective.
  - *Comment: the way BERT is trained is crucial for its final performance.*
- DistilBERT was trained on eight 16GB V100 GPUs for approximately three and a half days using the concatenation of Toronto Book Corpus and English Wikipedia (same data as original BERT).

# Recent achievements: GPT-3

GPT-3, May 2020  
175B parameters  
(proportions are incorrect for visual sake)

number of parameters, millions

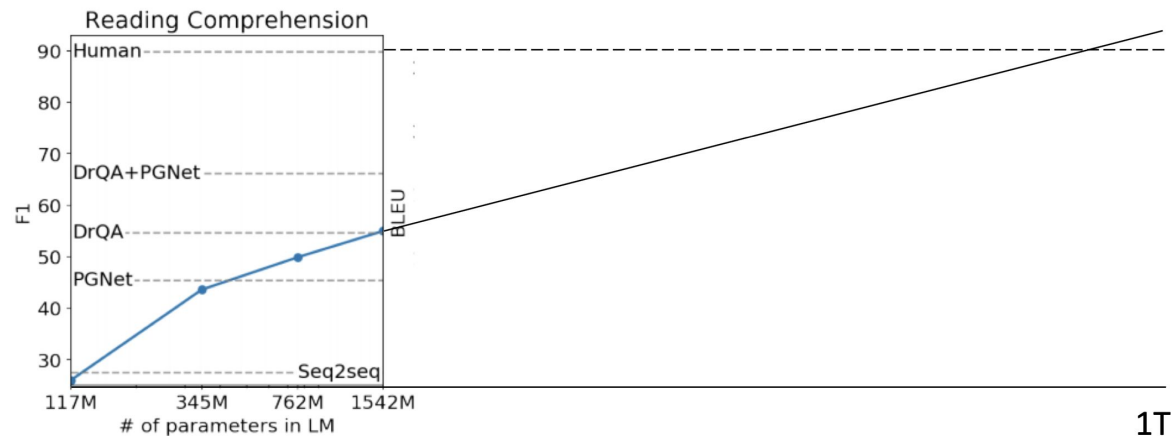


# Recent achievements: GPT-3

GPT-3, May 2020

175B parameters

(proportions are incorrect for visual sake)



Hypothesis from Stanford CS224n (2019) lecture 20



# Bonus: FusionNet

MLP (Additive) form:

$$S_{ij} = s^T \tanh(W_1 c_i + W_2 q_j)$$

Space:  $O(mnk)$ ,  $W$  is  $k \times d$

Bilinear (Product) form:

$$S_{ij} = c_i^T W q_j$$

$$S_{ij} = c_i^T U^T V q_j$$

Space:  $O((m+n)k)$

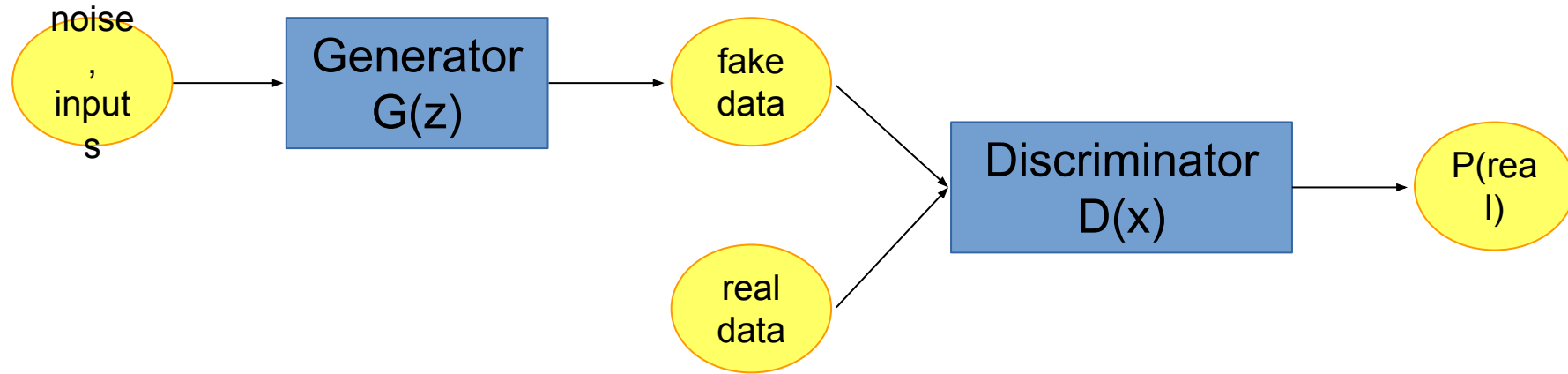
$$S_{ij} = c_i^T W^T D W q_j$$

1. Smaller space
2. Non-linearity

$$S_{ij} = \text{Relu}(c_i^T W^T) D \text{Relu}(W q_j)$$

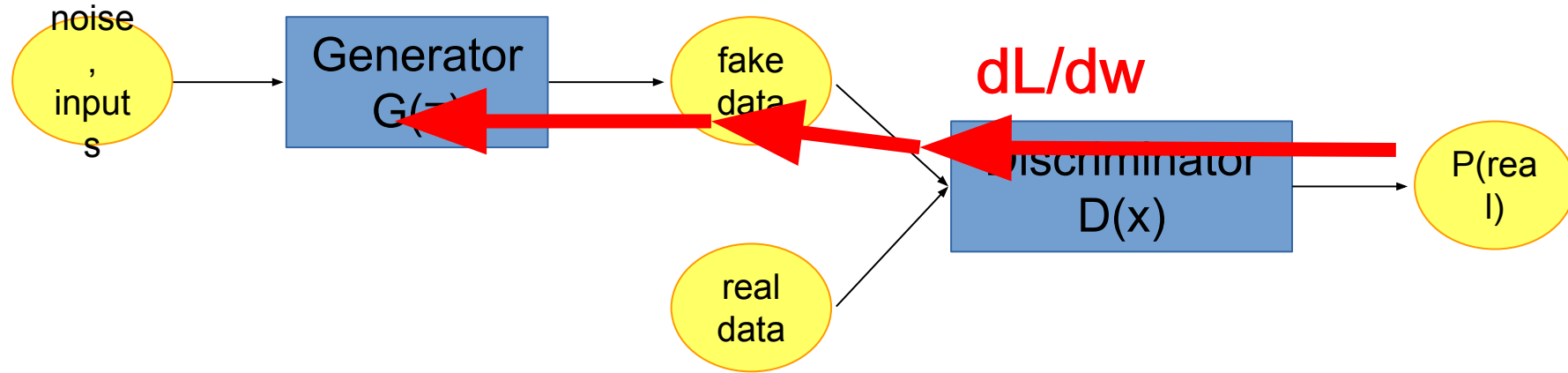


## Generalized GAN scheme



# Bonus: discrete GANs

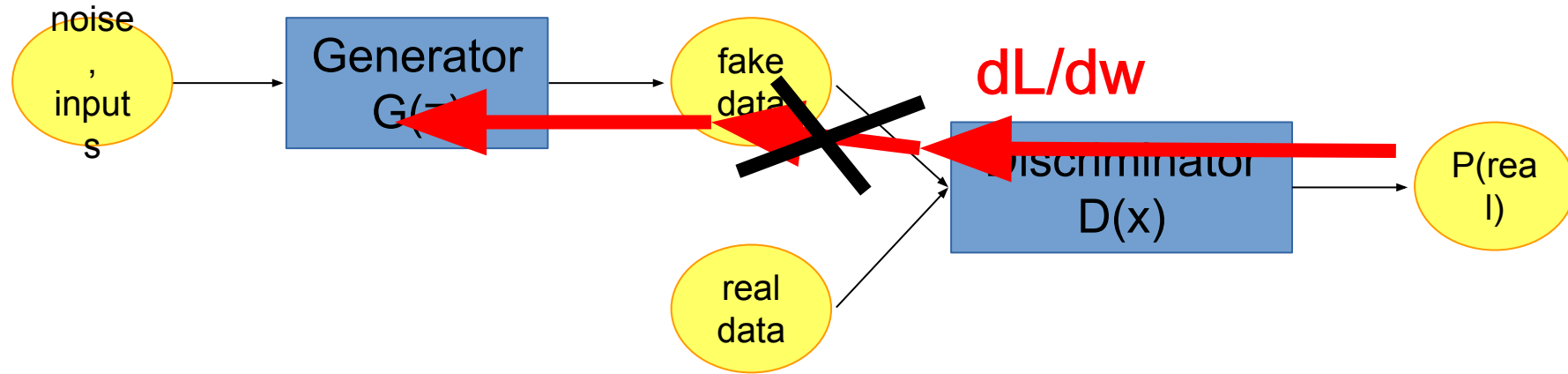
## Generalized GAN scheme



# Bonus: discrete GANs

Standard scheme fails if  $G(z)$  is discrete

- generating text
- generating music notes
- generating molecules
- binary image masks



We can train generator with Reinforcement Learning methods!

$$\nabla J = E_{\substack{z \sim p(z) \\ x \sim P(x|G_\theta(z))}} \nabla \log P(x|G_\theta(z)) D(x)$$