

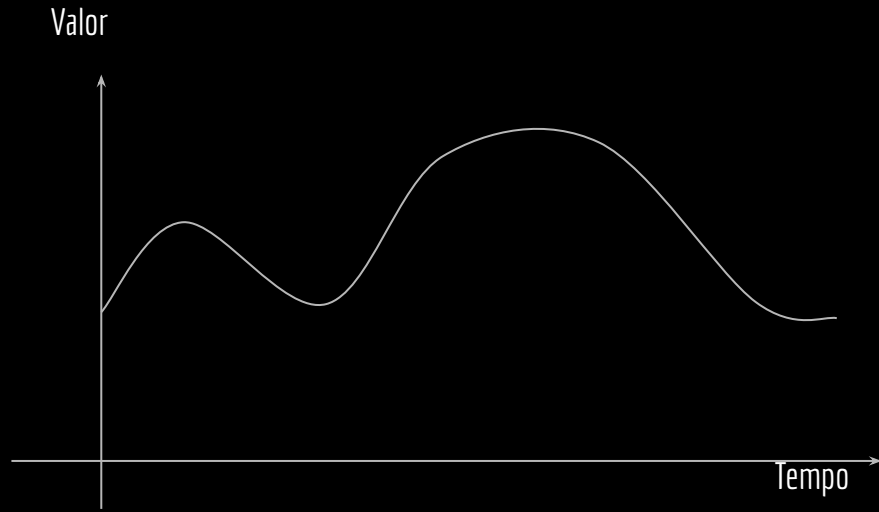
Alexa, faça de conta que você não está escutando.

Transformers Multimodais

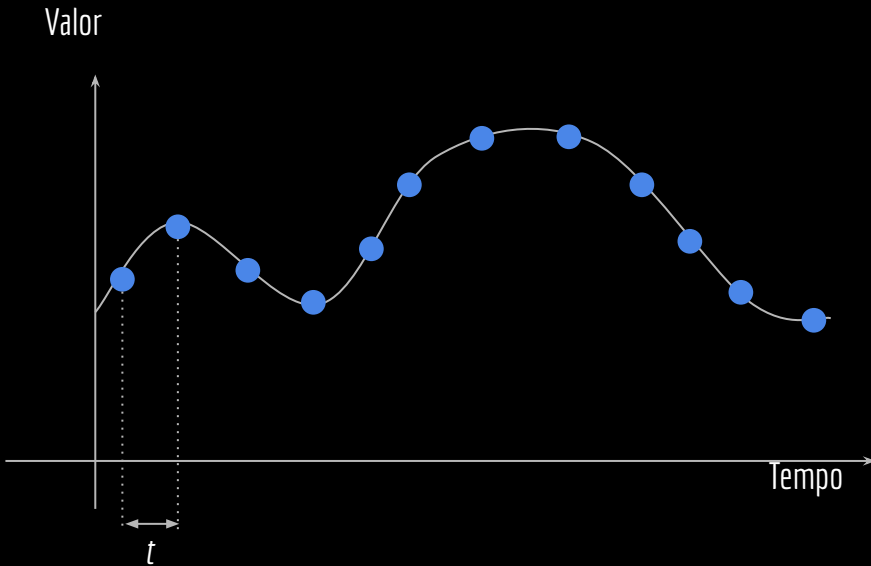
Paulo Ricardo Lisboa de Almeida



Mundo analógico e discreto



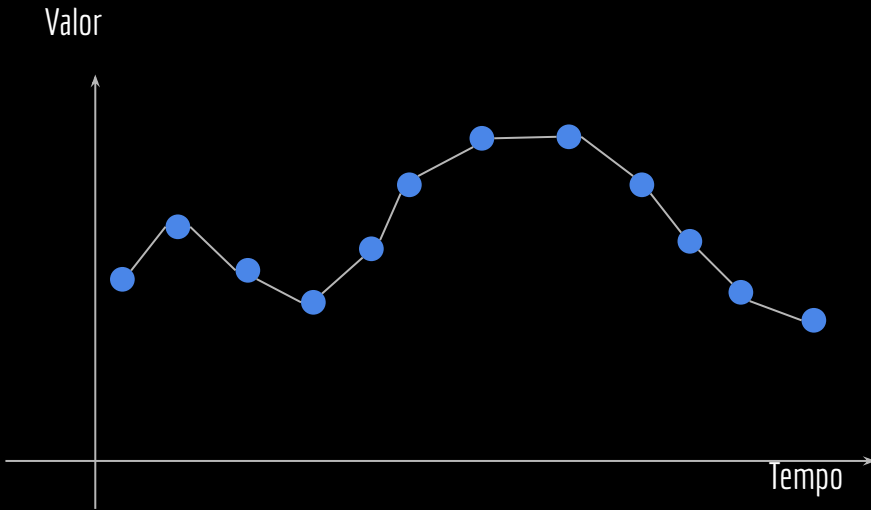
Mundo analógico e discreto



Amostragem.

Podemos discretizar um sinal analógico, verificando o seu valor a cada t instantes de tempo (t é o período).

Mundo analógico e discreto



Amostragem.

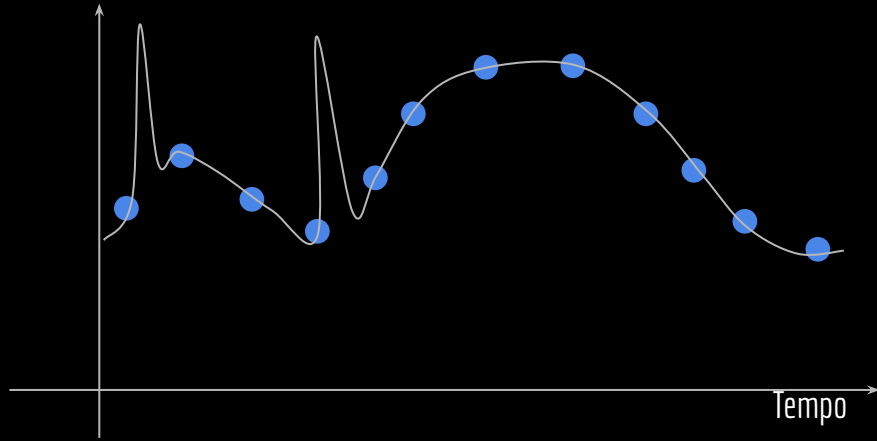
Podemos discretizar um sinal analógico, verificando o seu valor a cada t instantes de tempo.

Quanto menor for t , mais fiel é o sinal discreto quando comparado ao analógico.

Podemos interpolar o sinal discreto para reconstruir (algo próximo) do sinal original.

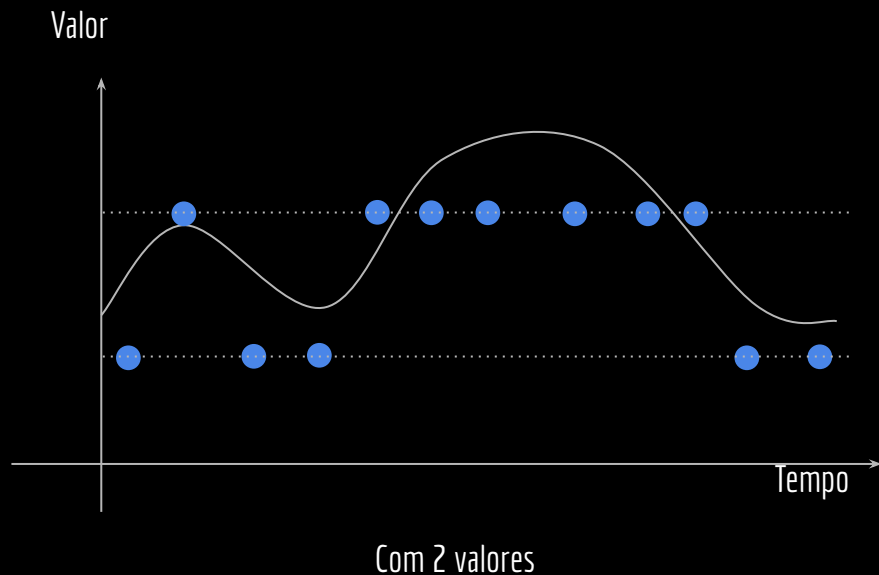
Mundo analógico e discreto

Valor



Durante a discretização, algumas informações podem ser perdidas.

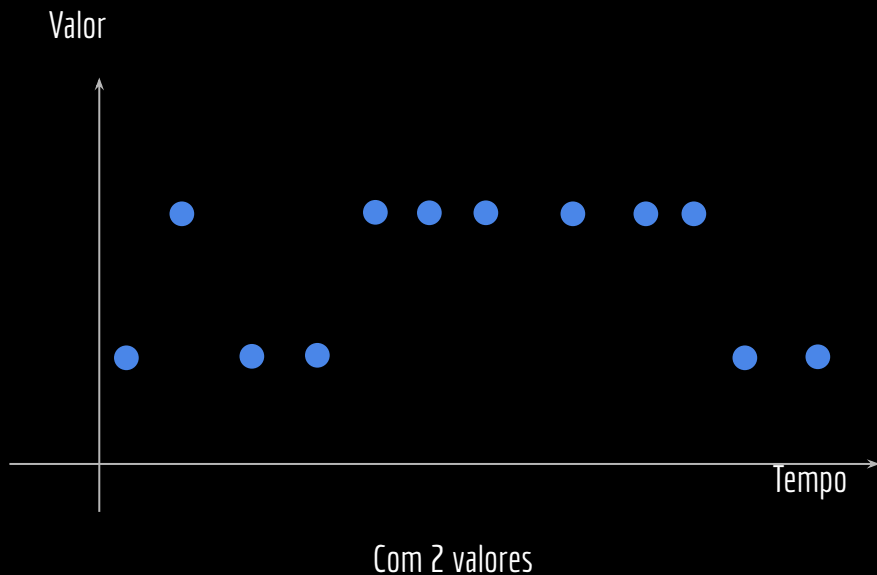
Mundo analógico e discreto



Quantização.

No eixo y, quantos valores podemos representar? Quanto mais valores valores, maior a fidelidade no eixo y com o sinal original, mas o hardware se torna mais complicado.

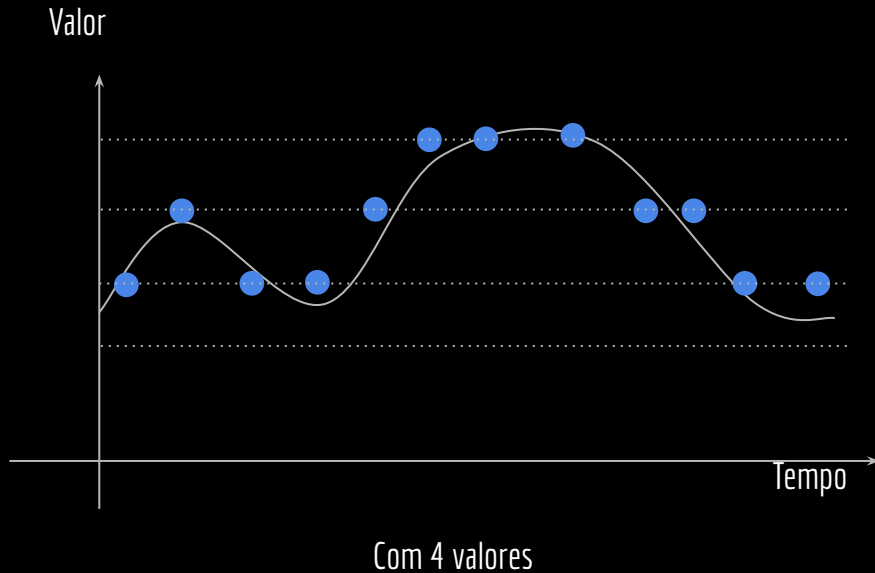
Mundo analógico e discreto



Quantização.

No eixo y, quantos valores podemos representar? Quanto mais valores valores, maior a fidelidade no eixo y com o sinal original, mas o hardware se torna mais complicado.

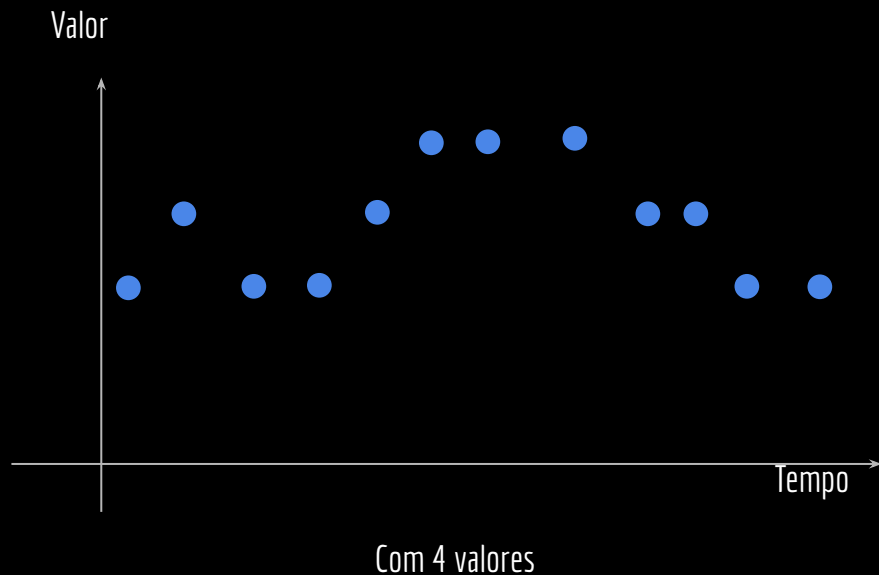
Mundo analógico e discreto



Quantização.

No eixo y, quantos valores podemos representar? Quanto mais valores valores, maior a fidelidade no eixo y com o sinal original, mas o hardware se torna mais complicado.

Mundo analógico e discreto



Quantização.

No eixo y, quantos valores podemos representar? Quanto mais valores valores, maior a fidelidade no eixo y com o sinal original, mas o hardware se torna mais complicado.

Exemplo

Geralmente as músicas que ouvimos possuem uma amostragem de 44.100Hz.

Isso significa que a cada $t = 1/441000 \approx 0,000023$ segundos um sinal é amostrado.

Exemplo

Geralmente as músicas que ouvimos possuem uma amostragem de 44.100Hz.

Isso significa que a cada $t = 1/441000 \approx 0,000023$ segundos um sinal é amostrado.

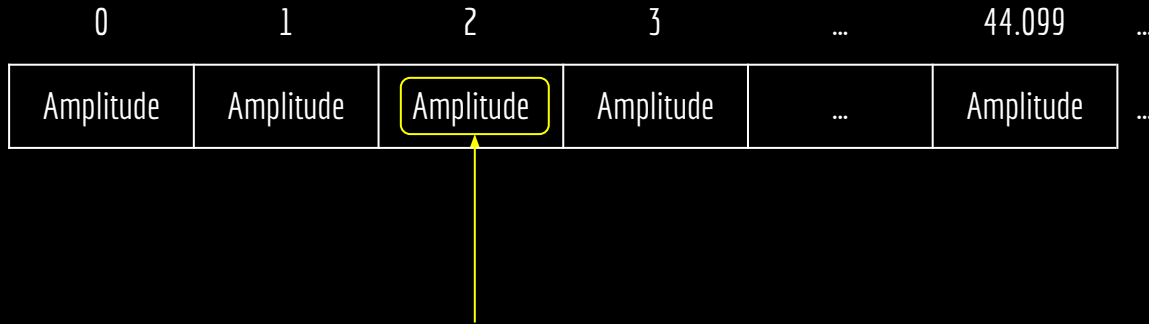
$f = 1/t$, onde f é a frequência, e t o período.

Exemplo - CD



Técnica chamada de Pulse-Code Modulation.

Exemplo - CD



Cada amostra de amplitude é armazenada como um valor de 16 bits.

Técnica chamada de Pulse-Code Modulation.

Exemplo - CD

Podemos adicionar mais fluxos para ter múltiplos canais de áudio. Por exemplo, em um CD com dois canais (stereo).

0	1	2	3	...	44.099	...
Amplitude	Amplitude	Amplitude	Amplitude	...	Amplitude	...

Canal 1.

0	1	2	3	...	44.099	...
Amplitude	Amplitude	Amplitude	Amplitude	...	Amplitude	...

Canal 2.

Técnica chamada de Pulse-Code Modulation.

Convoluções

0	1	2	3	...	44.099	...
Amplitude	Amplitude	Amplitude	Amplitude	...	Amplitude	...

Podemos aplicar uma CNN.

Convoluções de $I \times N \times C$, onde C é o número de canais, e N é o tamanho do Kernel.

Transformer

0	1	2	3	...	44.099	...
Amplitude	Amplitude	Amplitude	Amplitude	...	Amplitude	...

Ou podemos aplicar um transformer, que vai ser capaz de capturar o contexto global.
Mas possui um custo quadrático de acordo com o tamanho da entrada.

Ideia

A mesma ideia aplicada para modelos de linguagem pode ser aplicada.

Dado um conjunto de tokens (de áudio) prever, por exemplo, qual o token (de texto) pertence a esse trecho.

Nesse exemplo, o modelo se torna um *speech to text*.

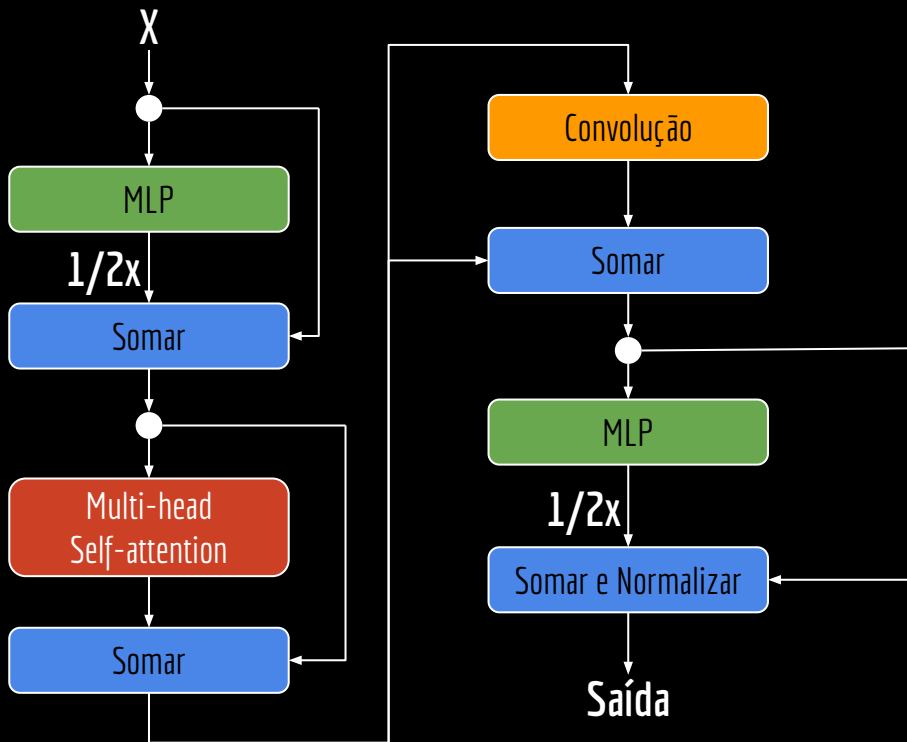
Estado da arte

Técnicas do estado da arte combinam convoluções com transformers.

Convoluções para capturar as correlações locais e reduzir a dimensionalidade dos dados.

Transformers para capturar as correlações globais.

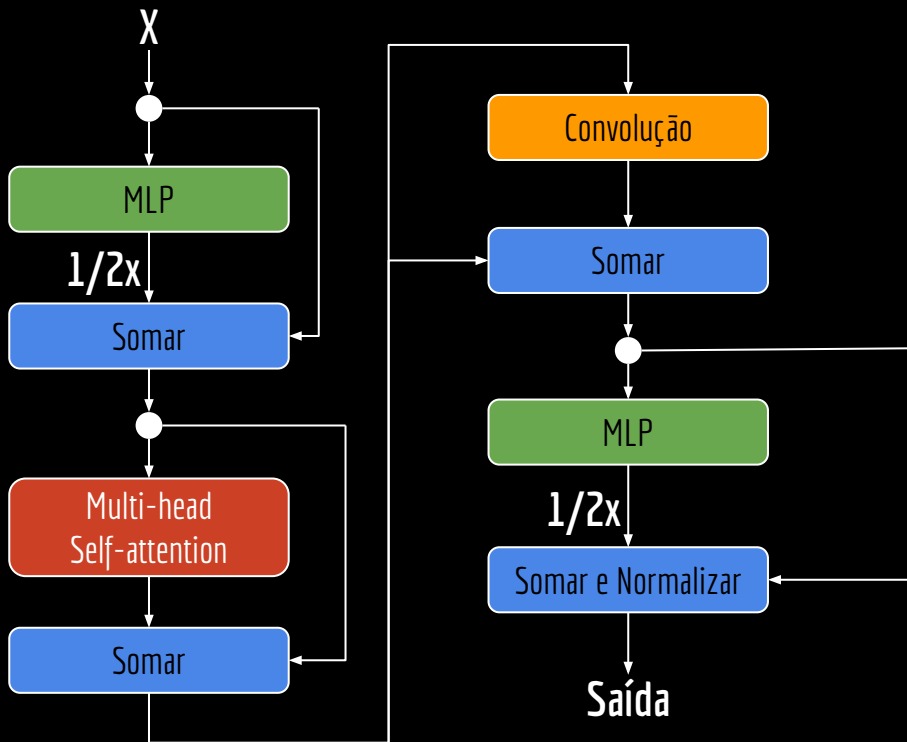
Conformer



GULATI, Anmol et al. Conformer: Convolution-augmented transformer for speech recognition. 2020.



Conformer



GULATI, Anmol et al. Conformer: Convolution-augmented transformer for speech recognition. 2020.

$$\tilde{x} = x + \frac{1}{2}MLP(x)$$

$$x' = \tilde{x} + MHSA(\tilde{x})$$

$$x'' = x' + Conv(x')$$

$$y = \text{normalizar}(x'' + \frac{1}{2}MLP(x''))$$



Conformer

Cabeças de atenção ficam em um “sanduíche” de MLPs.

O primeiro MLP projeta os dados para uma dimensão 4x menor. O segundo desfaz essa operação.

GULATI, Anmol et al. Conformer: Convolution-augmented transformer for speech recognition. 2020.



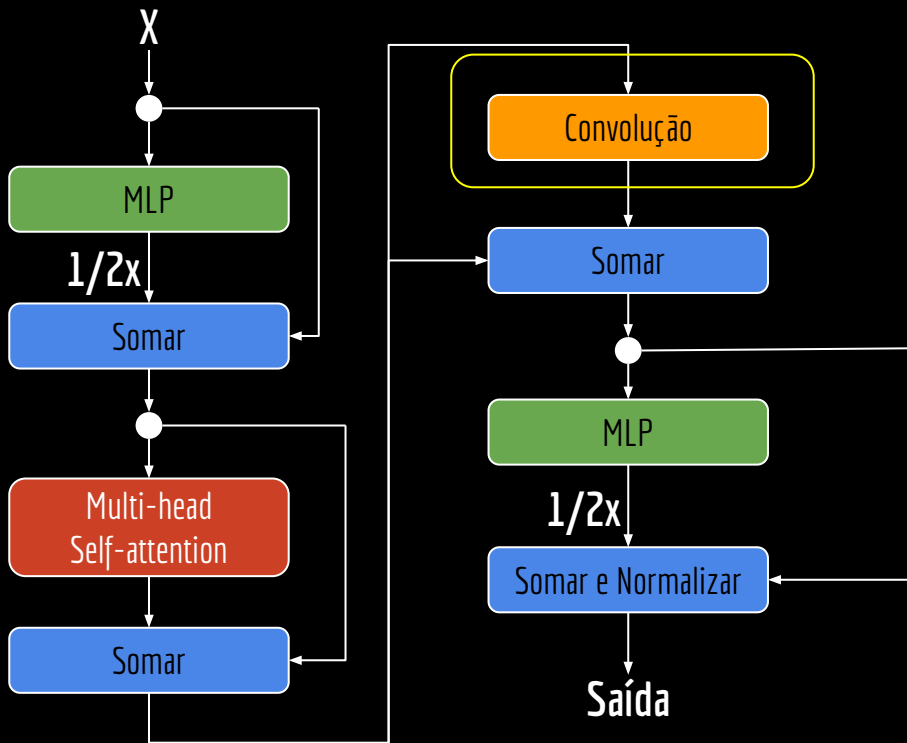
$$\tilde{\mathbf{x}} = \mathbf{x} + \frac{1}{2}MLP(\mathbf{x})$$

$$\mathbf{x}' = \tilde{\mathbf{x}} + MHSA(\tilde{\mathbf{x}})$$

$$\mathbf{x}'' = \mathbf{x}' + Conv(\mathbf{x}')$$

$$\mathbf{y} = \text{normalizar}(\mathbf{x}'' + \frac{1}{2}MLP(\mathbf{x}''))$$

Conformer



GULATI, Anmol et al. Conformer: Convolution-augmented transformer for speech recognition. 2020.



Camada convolucional aplica convolução de tamanho 1 (pointwise) para combinar os canais, e depois aplica uma convolução de tamanho 31 para capturar as correlações locais.

Projeto

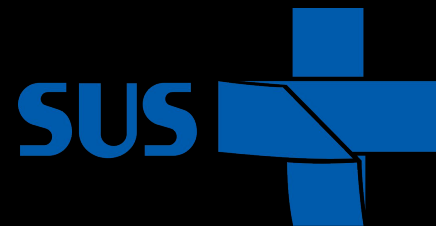
O uso de transformers está sendo estudado no Projeto “Comunicação direta entre usuários e a Atenção Primária à Saúde (APS)”.

Pesquisa atual: como converter as respostas de usuários dadas via ligação de voz para dados textuais.

Algo parecido com um “Quiz” que a pessoa pode responder utilizando a própria voz, sem formulários.

Os dados podem ser processados para identificar possíveis melhorias em determinadas regiões, hospitais, postos de saúde, ...

Obs.: os dados são processados apenas com o consentimento dos usuários.



Outras modalidades

Dados que podemos projetar em um espaço (*embedding*) de tokens, e onde há uma dependência sequencial entre os dados, podem ser processados através de transformers.

O processamento de áudio foi um exemplo.

Outras modalidades

Dados que podemos projetar em um espaço (*embedding*) de tokens, e onde há uma dependência sequencial entre os dados, podem ser processados através de transformers.

O processamento de áudio foi um exemplo.

Outro exemplo: processamento de imagens.

Transformers para visão

Em visão, podemos utilizar cada pixel individual como um *token*.

Vai ser computacionalmente custoso, mas funciona.

Transformers para visão

Em visão, podemos utilizar cada pixel individual como um *token*.

Vai ser computacionalmente custoso, mas funciona.

Para reduzir o custo podemos, por exemplo:

- Utilizar patches de $N \times N$ das imagens, e projetar cada patch como um token, ou
- Utilizar filtros de convolução para reduzir a dimensionalidade das imagens originais.

Ideia do ViT: *Vision Transformer*

Transformers para visão

O embedding final gerado pelos transformers pode ser usado para, por exemplo, passar por um MLP para classificar a imagem.

Ou podemos montar um modelo baseado em auto regressão:

Dado o início de uma imagem, prever os próximos pixels dela (terminar de desenhar).

Multimodalidade

Ao transformar as entradas e saídas em tokens, é relativamente trivial unir múltiplos modais de informação em uma única rede.

Exemplos:

- Uma rede que recebe ao um prompt de texto, e uma imagem de base, e a modifica de acordo com o prompt.
- Uma rede que recebe um prompt, e gera uma imagem de saída (na forma de auto regressão, na próxima iteração, recebe o prompt, e o primeiro pixel já gerado, ...).

Exemplo - CM3Leon

YU, Lili et al. Scaling autoregressive multi-modal models: Pretraining and instruction tuning. arXiv preprint arXiv:2309.02591. 2023.

arXiv:2309.02591v1 [cs.LG] 8 Sep 2023

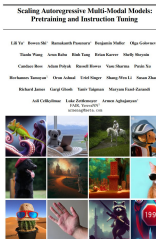


Figure 1: Results of CM3L on various multimodal generation tasks. Rows 1-2 show text-to-image generation, Row 3 shows image-to-image generation, Row 4 shows image-to-video generation, Row 5 shows image-to-audio generation, Row 6 shows image-to-text generation, Row 7 shows image-to-text-audio generation, Row 8 shows image-to-text-audio-visual generation.



Input



"What would she look like as a bearded man?"



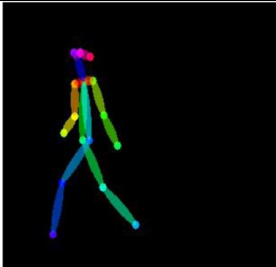
"Put on a pair of sunglasses"



"she should look 100 years old"



"Apply face paint"



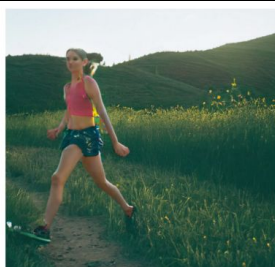
Extracted (openpose) pose



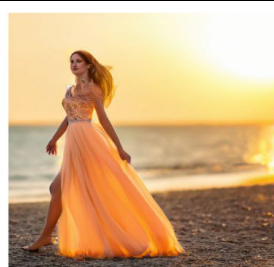
"Businessman in city street"



"A boy running on the grass of a soccer field"



"Young girl running on mountain trail with wild flowers"



"Beautiful women walking on the beach at sunset"

Timeline



1943 - Warren McCulloch e Walter Pitts.
Neurônio humano modelado como um
perceptron eletrônico.

Timeline

1943

Neurônio
Eletrônico



1957 - Frank Rosenblatt.
Implementação de um perceptron em um
computador IBM 704.

Timeline

1943

Neurônio
Eletrônico



1957

Perceptron
de
Rosenblatt



1970/1971 - Seppo Linnainmaa e Paul Werbos.
Desenvolvimento da retropropagação.

Timeline

1969 - Kunihiko Fukushima propõe a ativação ReLU.

1943

Neurônio Eletrônico

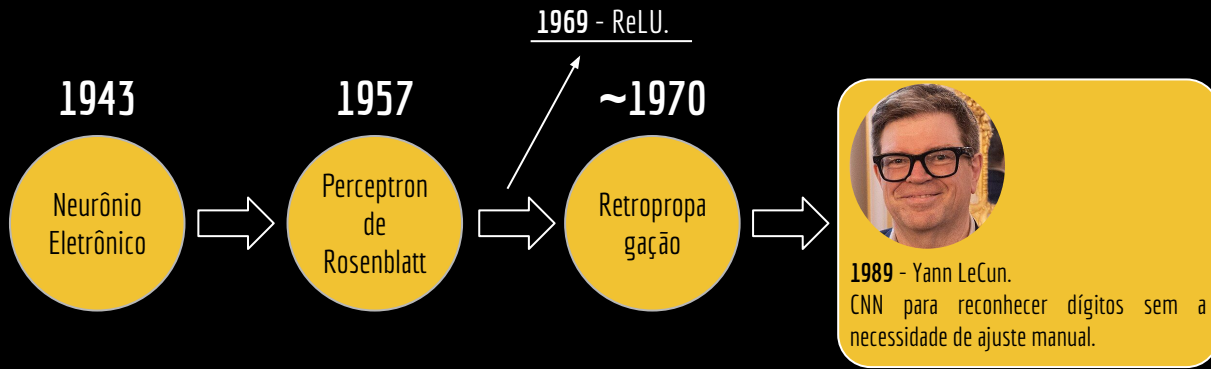
1957

Perceptron de Rosenblatt

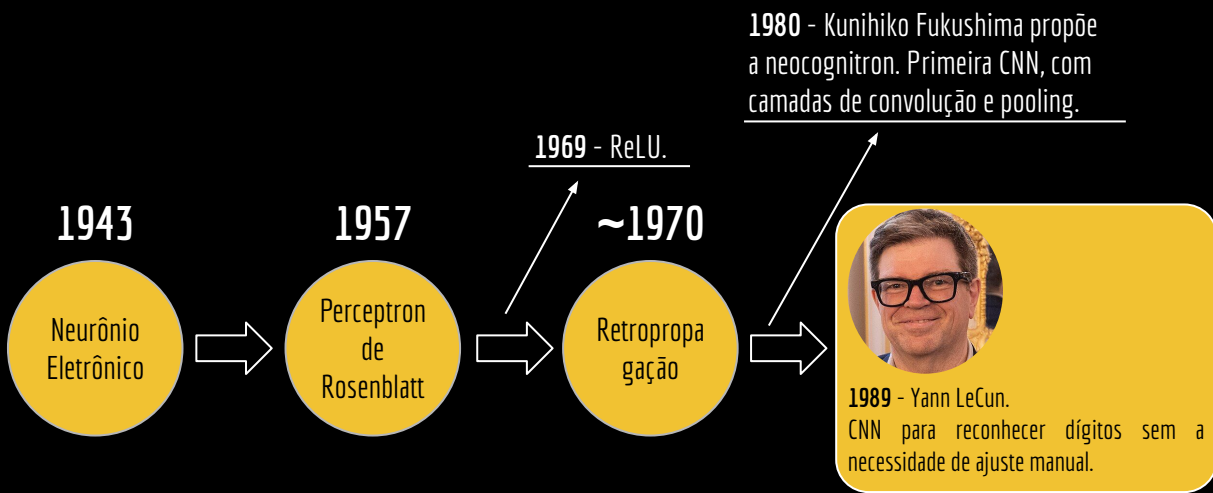


1970/1971 - Seppo Linnainmaa e Paul Werbos. Desenvolvimento da retropropagação.

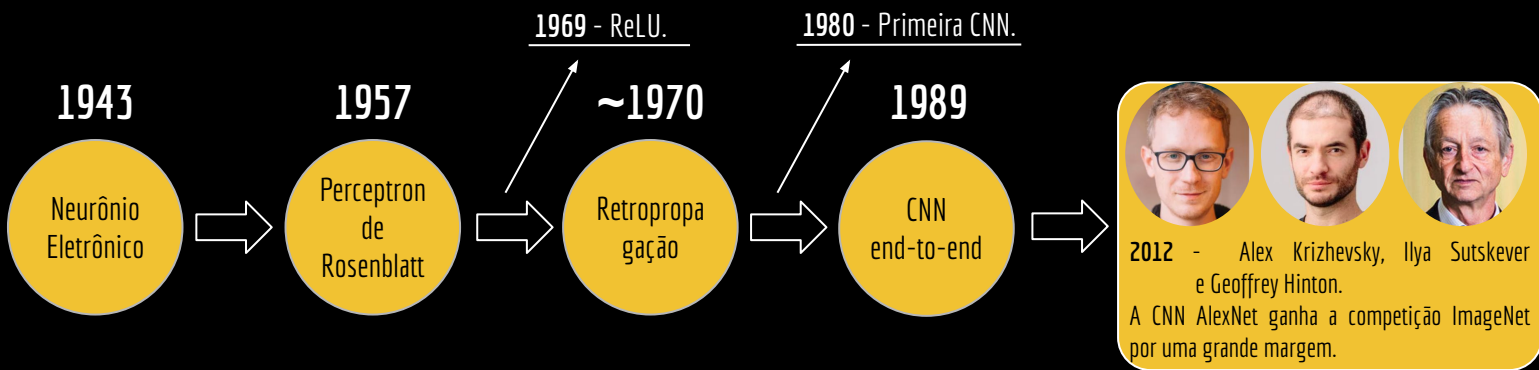
Timeline



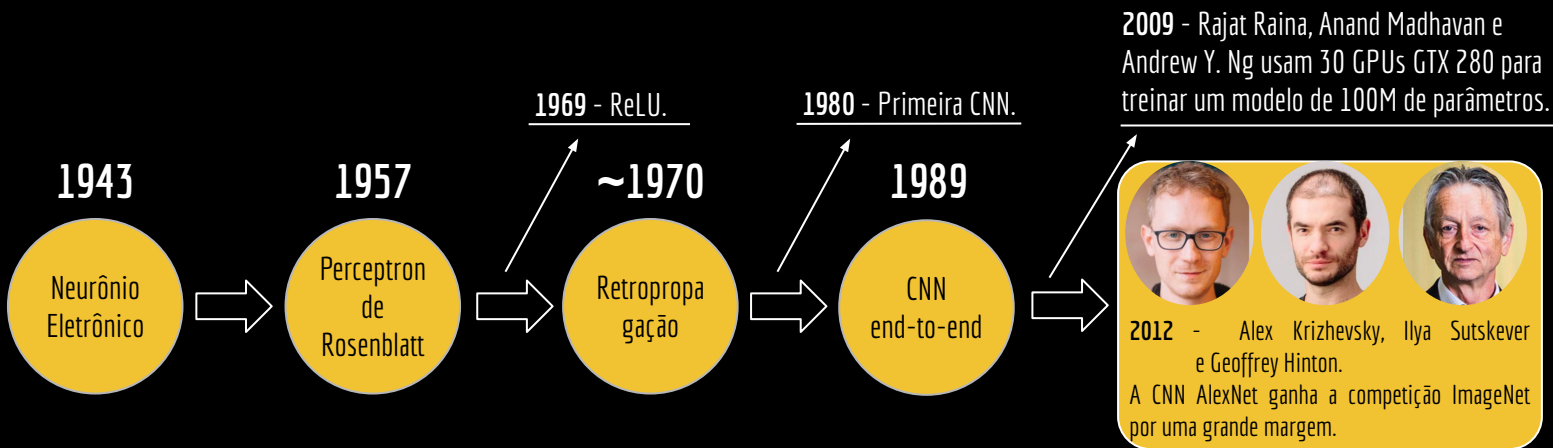
Timeline



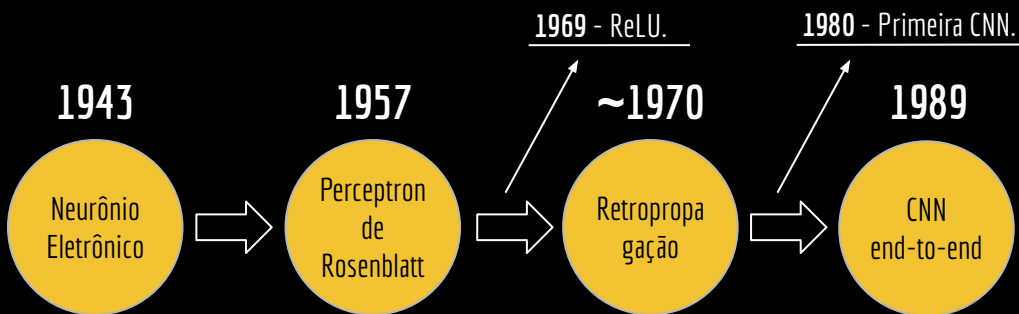
Timeline



Timeline



Timeline

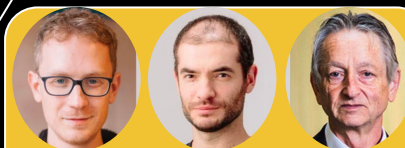


1969 - ReLU.

1980 - Primeira CNN.

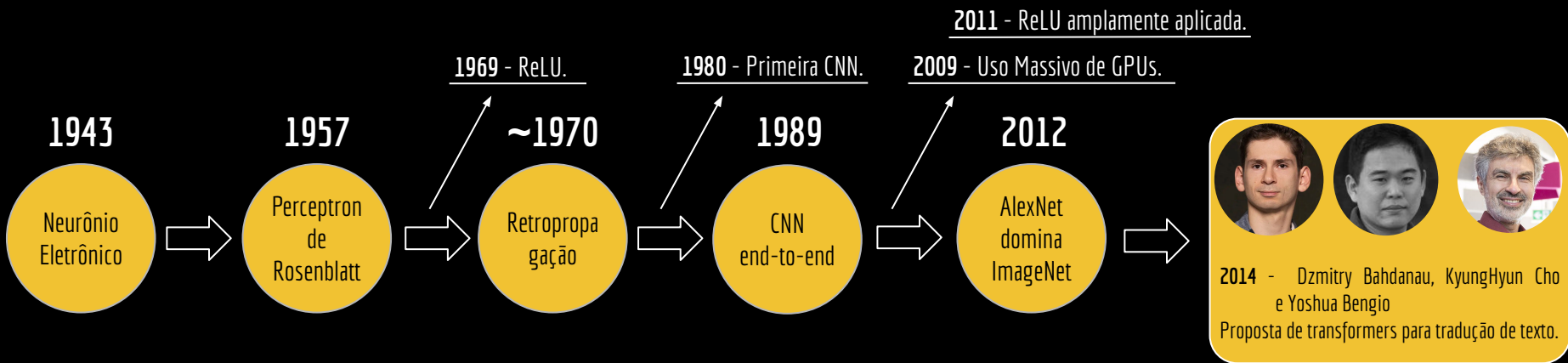
2011 - Xavier Glorot, Antoine Bordes e Yoshua Bengio descobrem que a ReLU funciona bem em redes profundas.

2009 - Rajat Raina, Anand Madhavan e Andrew Y. Ng usam 30 GPUs GTX 280 para treinar um modelo de 100M de parâmetros.

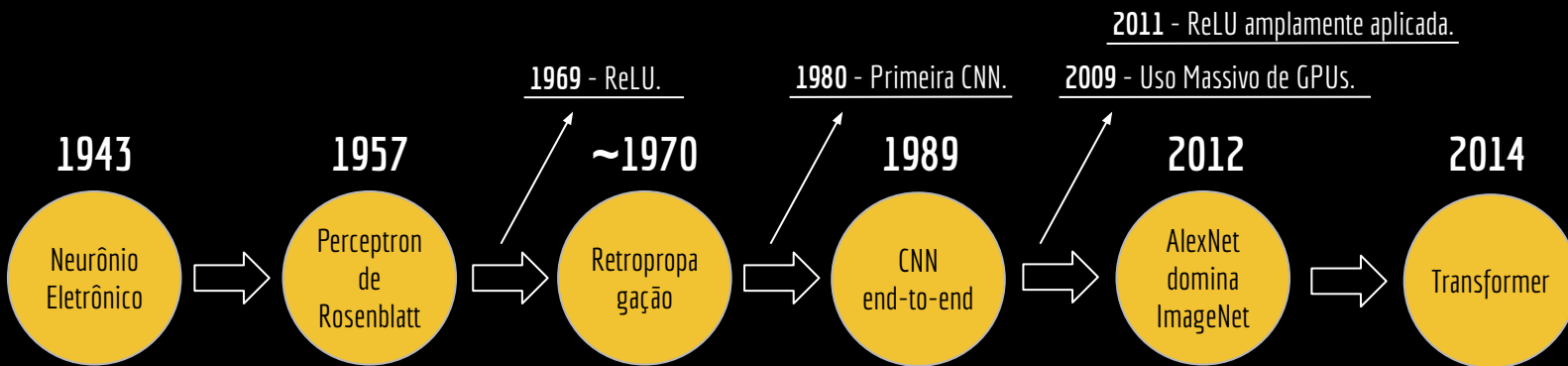


2012 - Alex Krizhevsky, Ilya Sutskever e Geoffrey Hinton.
A CNN AlexNet ganha a competição ImageNet por uma grande margem.

Timeline

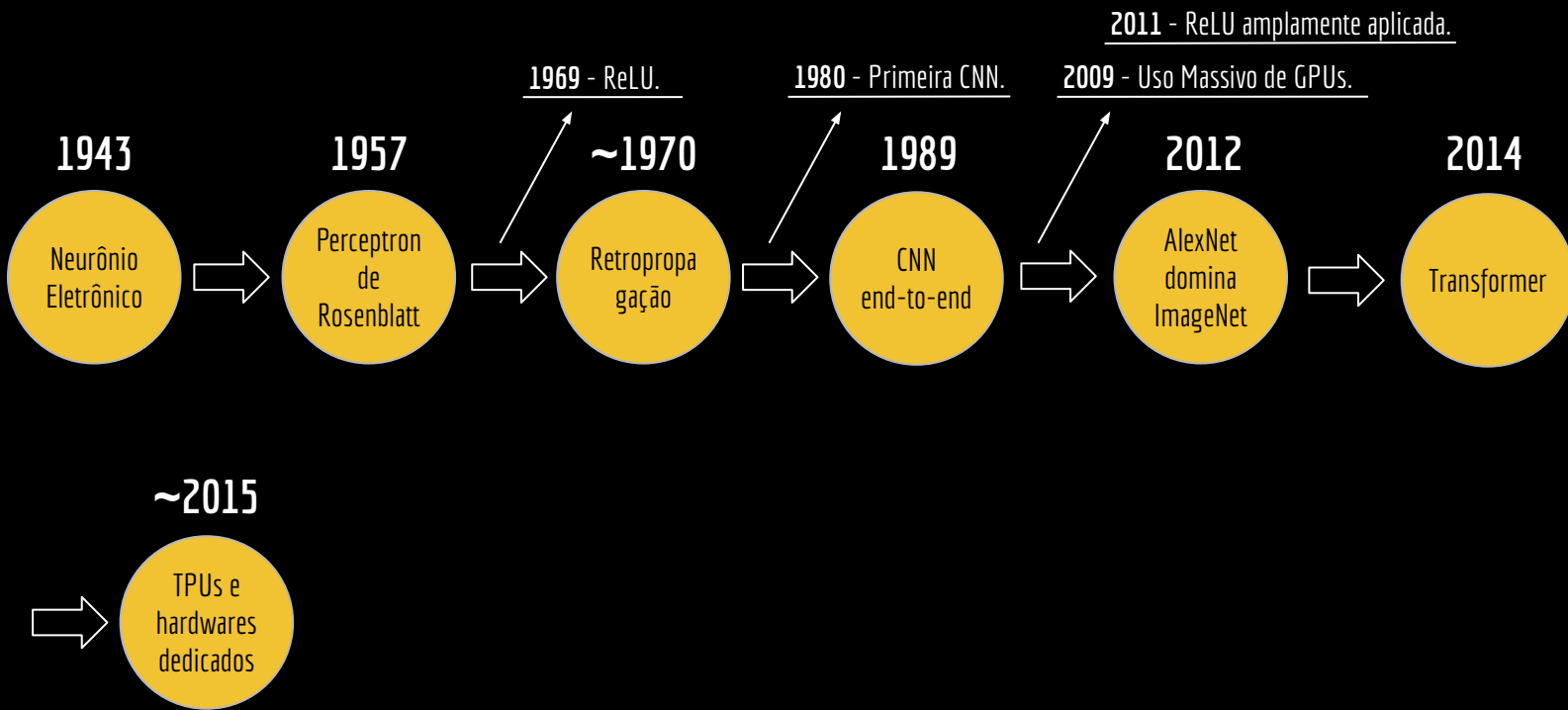


Timeline

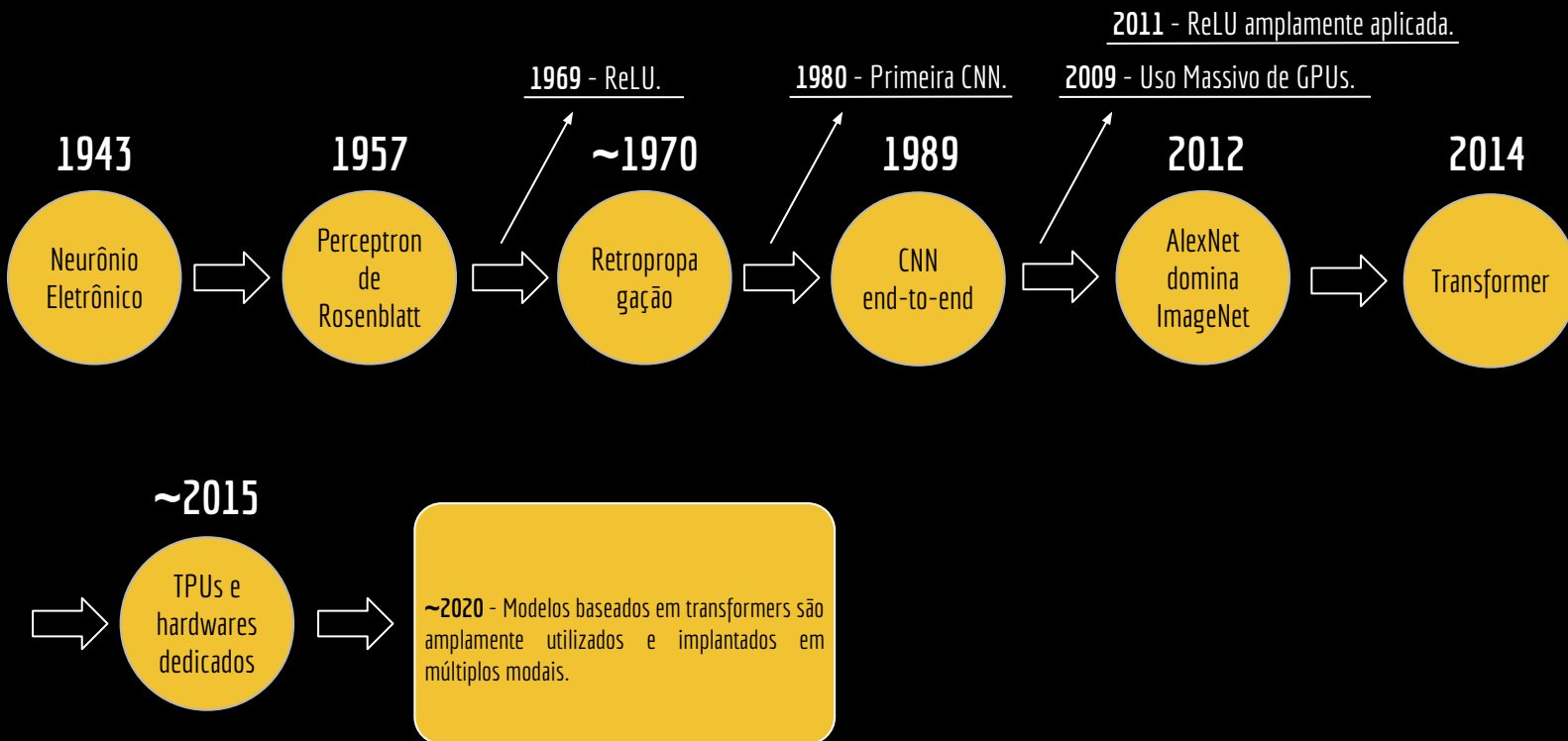


2015 - Patterson et. al.
Desenvolvimento de hardwares massivamente paralelos para ML, como TPUs.

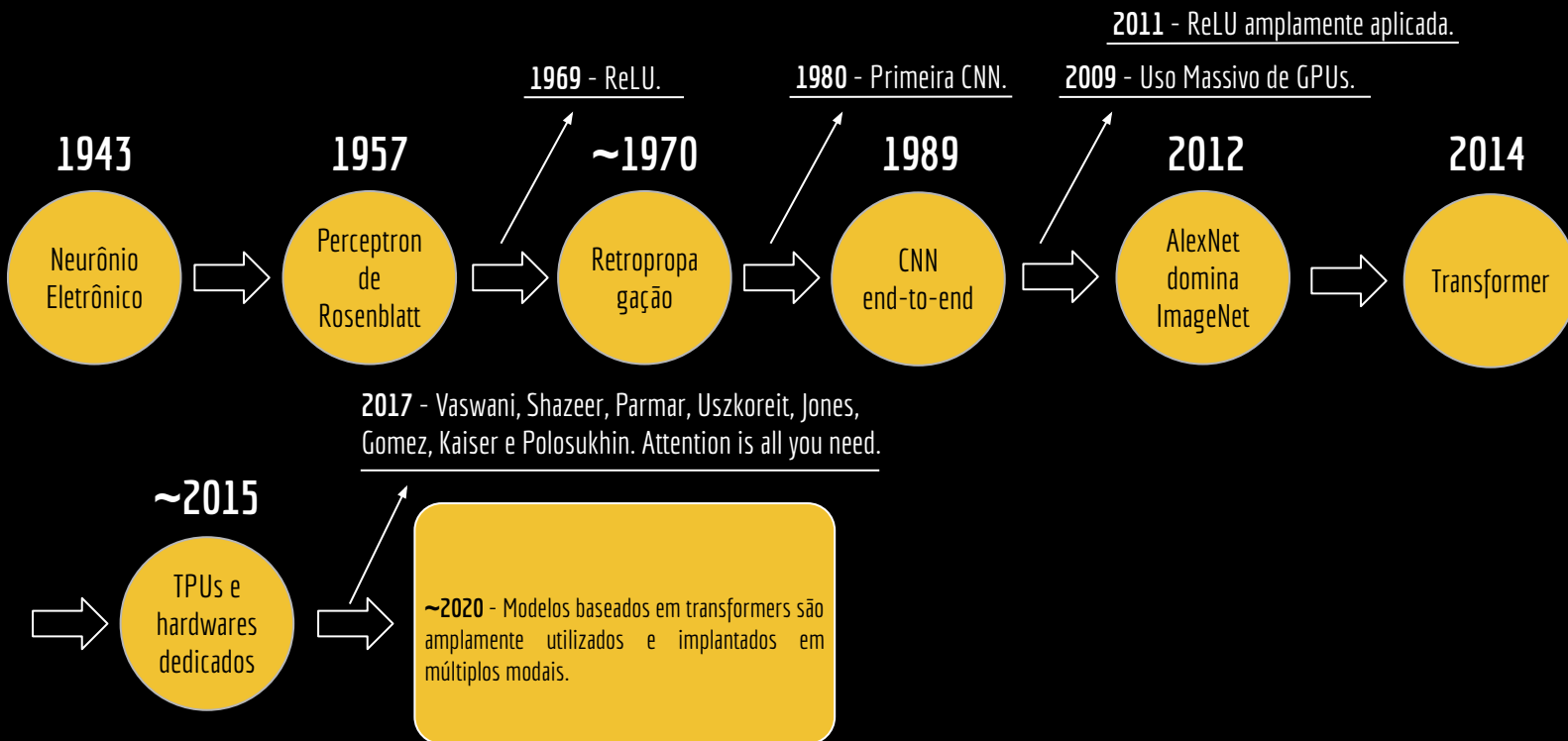
Timeline



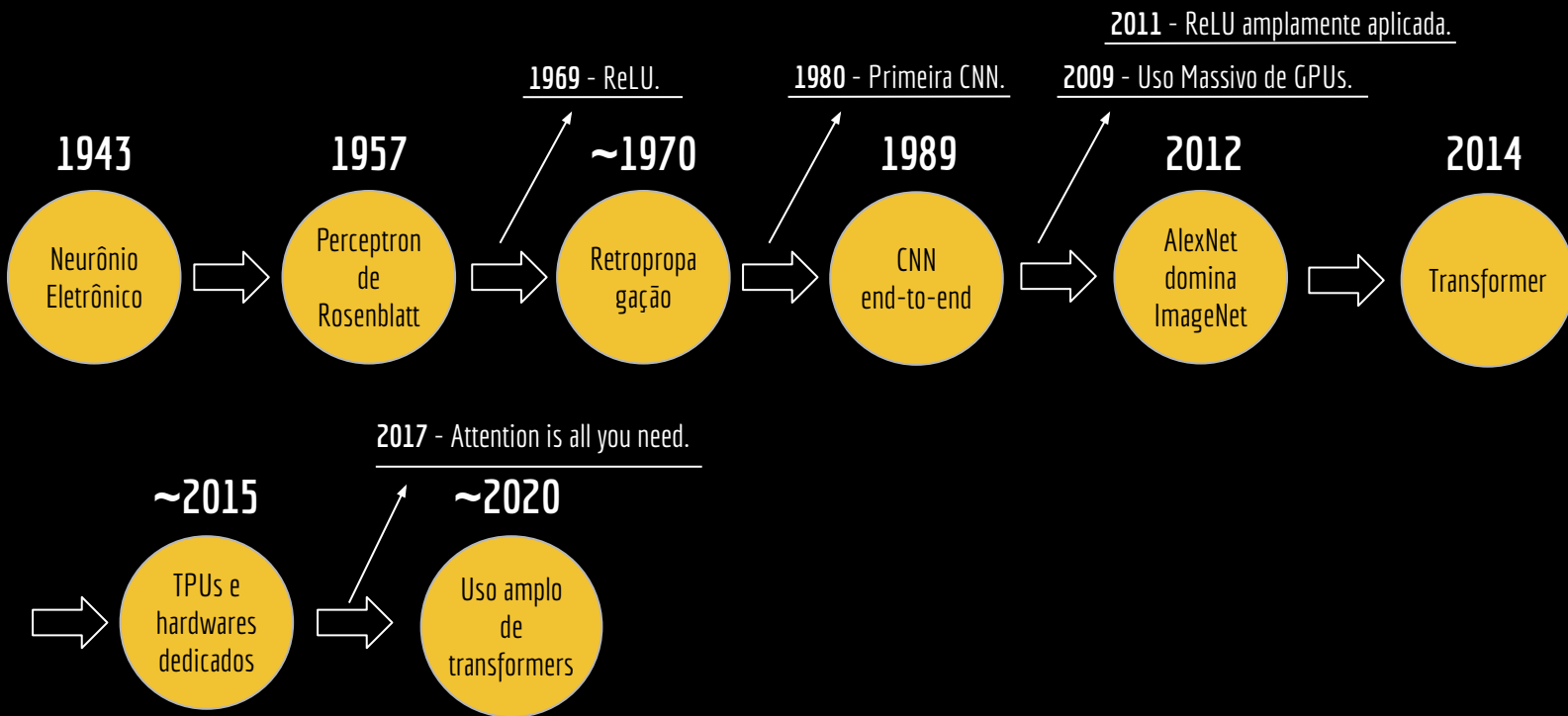
Timeline



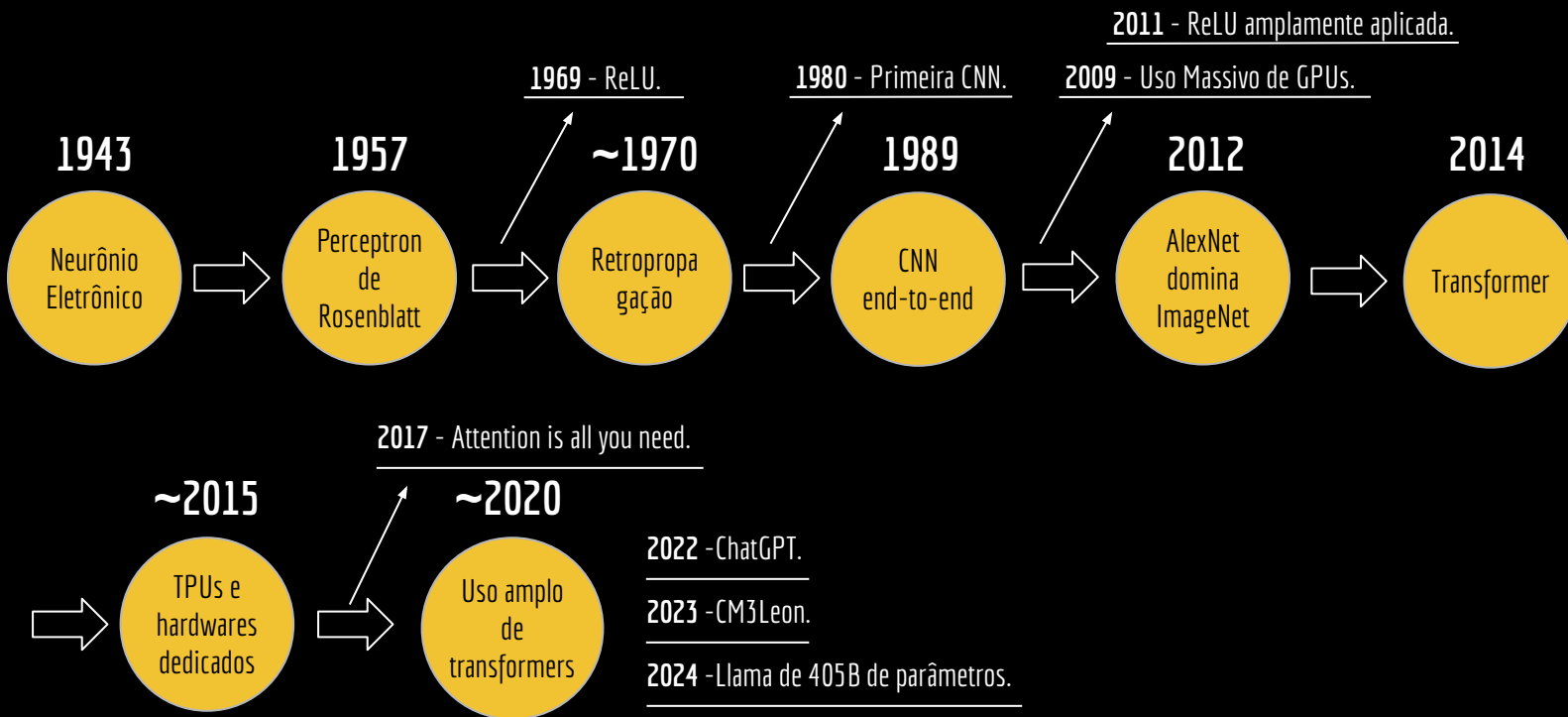
Timeline



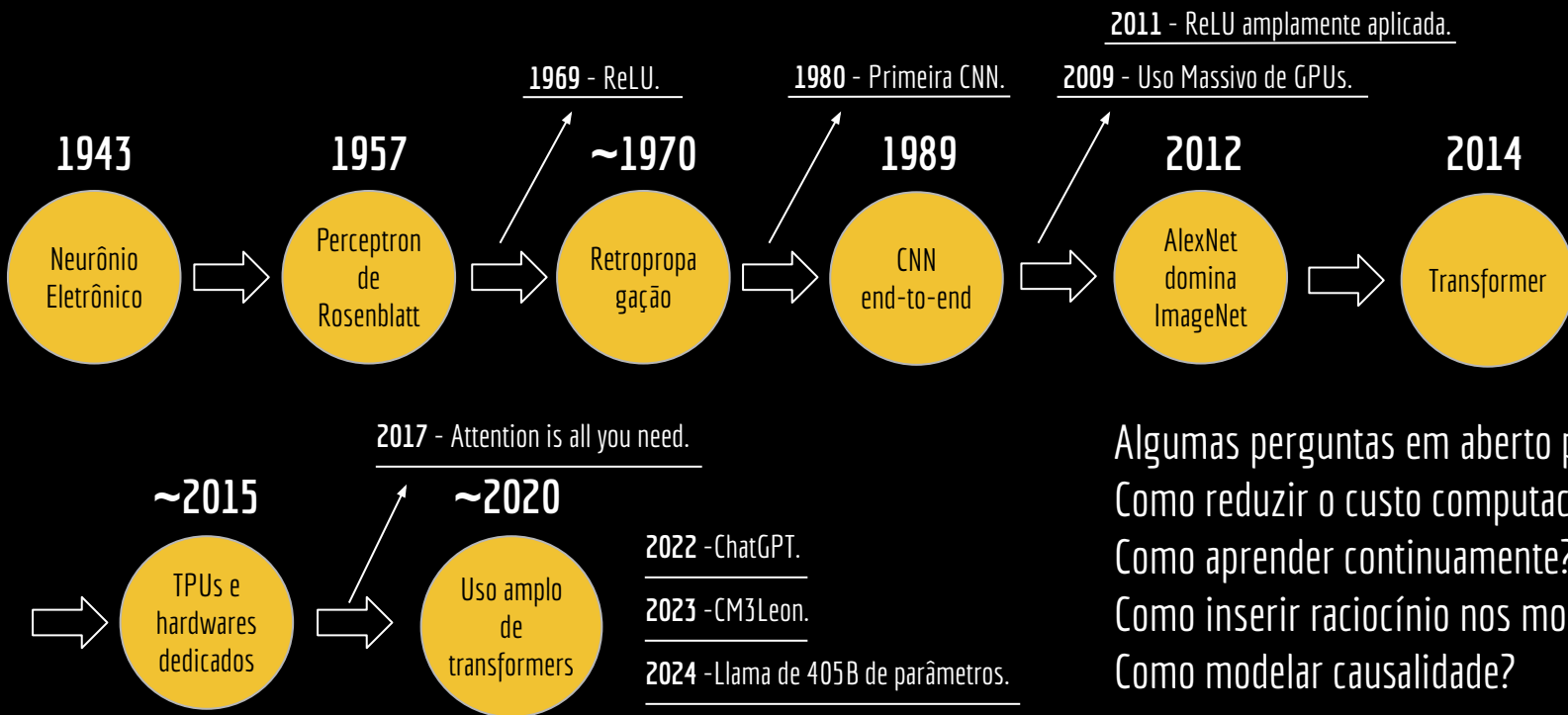
Timeline



Timeline



"That's all Folks!"



Algumas perguntas em aberto para > 2024:
Como reduzir o custo computacional?
Como aprender continuamente?
Como inserir raciocínio nos modelos?
Como modelar causalidade?

...

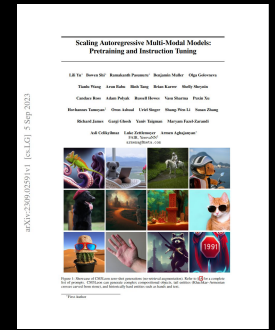
Referências

GULATI, Anmol et al. Conformer: Convolution-augmented transformer for speech recognition. 2020.

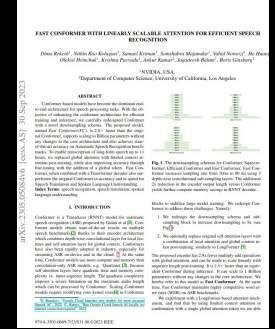


YAO, Zengwei et al. Zipformer: A faster and better encoder for automatic speech recognition. In: International Conference on Learning Representations. 2024.

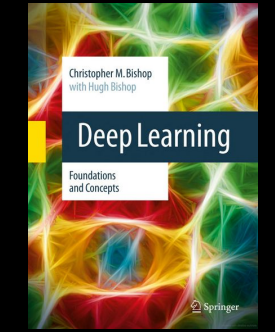
YU, Lili et al. Scaling autoregressive multi-modal models: Pretraining and instruction tuning. arXiv preprint arXiv:2309.02591. 2023.



REKESH, Dima et al. Fast conformer with linearly scalable attention for efficient speech recognition. IEEE Automatic Speech Recognition and Understanding Workshop. 2023.



Bishop, C. M., Bishop, H. Deep Learning: Foundations and Concepts. 2023.



Licença

Esta obra está licenciada com uma Licença [Creative Commons Atribuição 4.0 Internacional](https://creativecommons.org/licenses/by/4.0/).