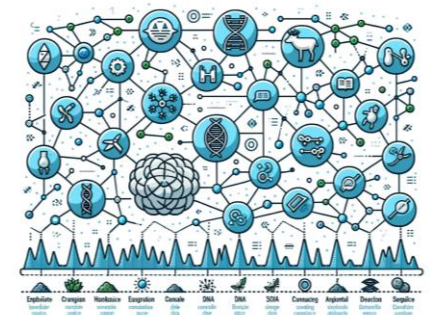


<https://github.com/Phillips-Lab-MTSU/DL-Workshop>

# EXPLORING DEEP LEARNING FOR SCIENTIFIC DISCOVERY

# Joshua L. Phillips

Associate Professor  
Department of Computer Science  
Program in Data Science  
Program in Computational and Data Science  
Middle Tennessee State University



ChatGPT 4 (DALL-E 3) Prompt: "Generate an image of a neural network processing multiple short strands of DNA from microbes taken from animal and soil samples." (2023-10-18)



ML algorithms need three things for learning:

1. Task
2. Performance Measure
3. Experience

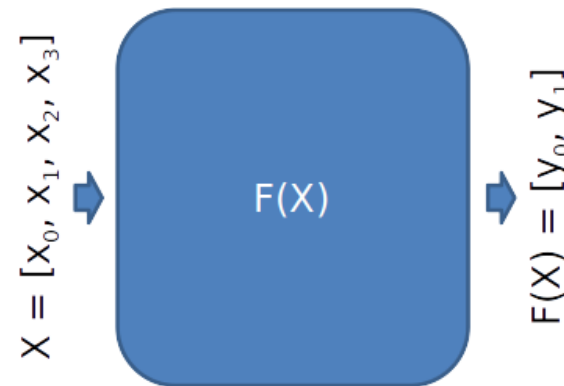
## MACHINE LEARNING BASICS

Q: How do we write computer programs that can solve problems which are very difficult?

A: We don't. We write computer programs that *learn* to solve problems for which it is very difficult to write computer programs.

### Learning Paradigms

- Unsupervised Learning
- Supervised Learning
- Reinforcement Learning
- Adversarial Learning

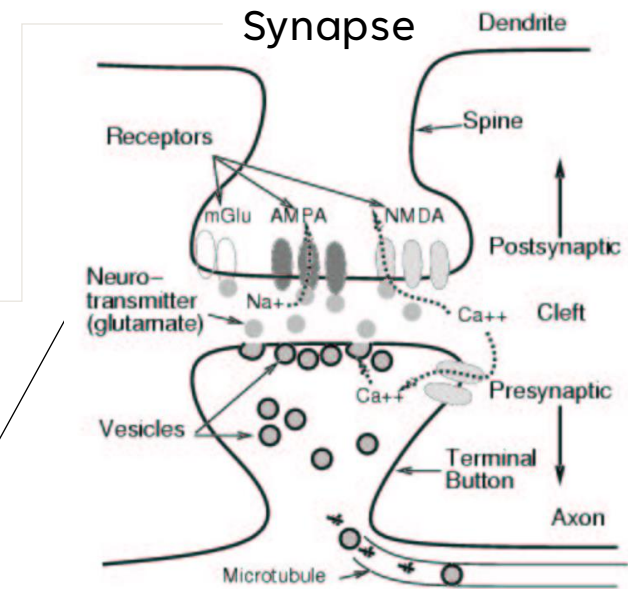
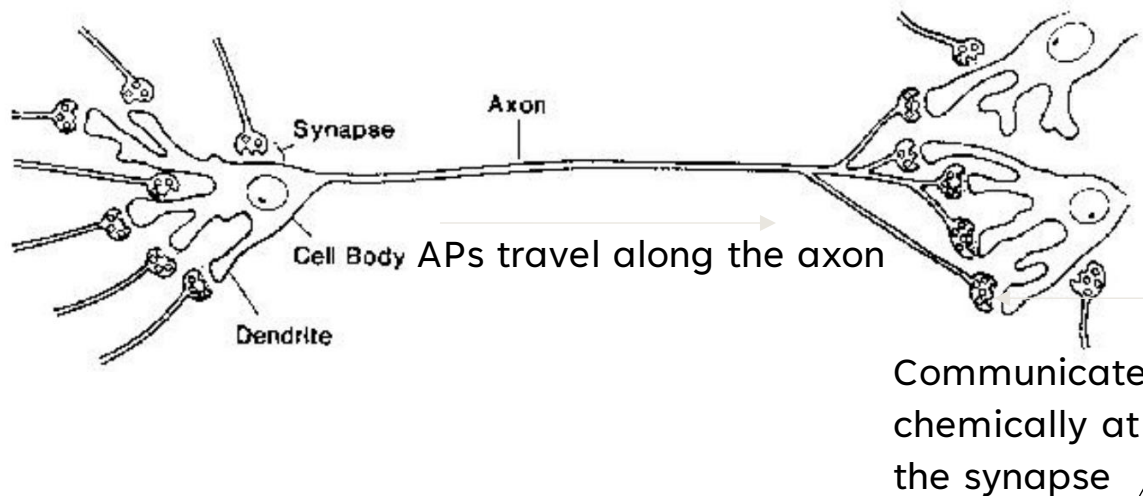
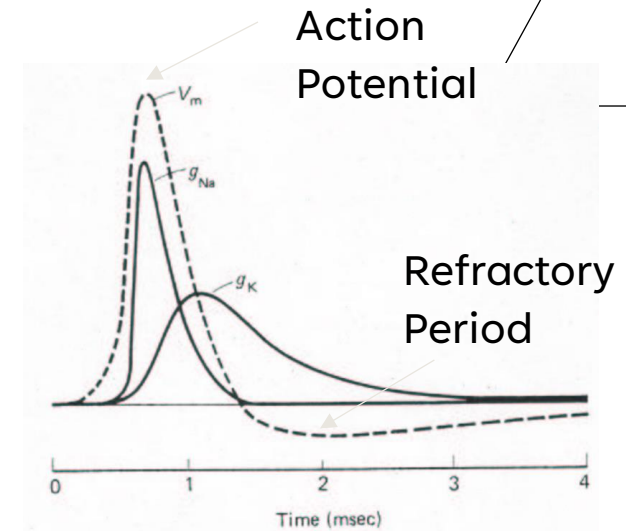


Deep Learning uses an architecture built of artificial neurons – if the neurons are wired together correctly, task is solved



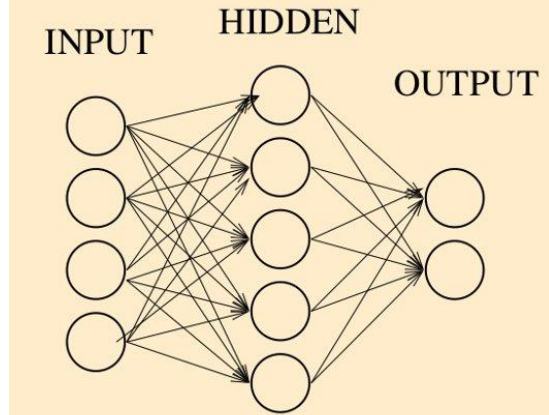
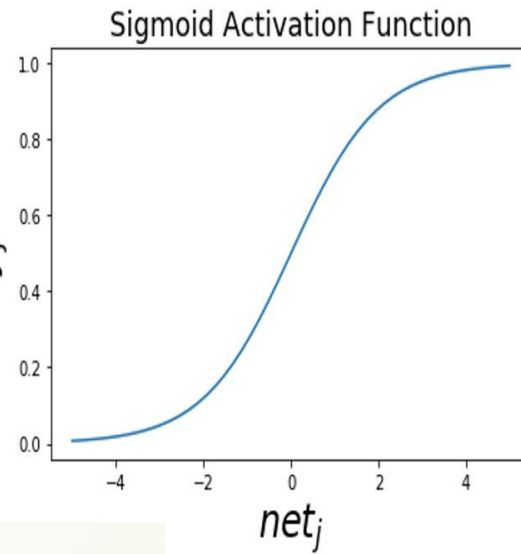
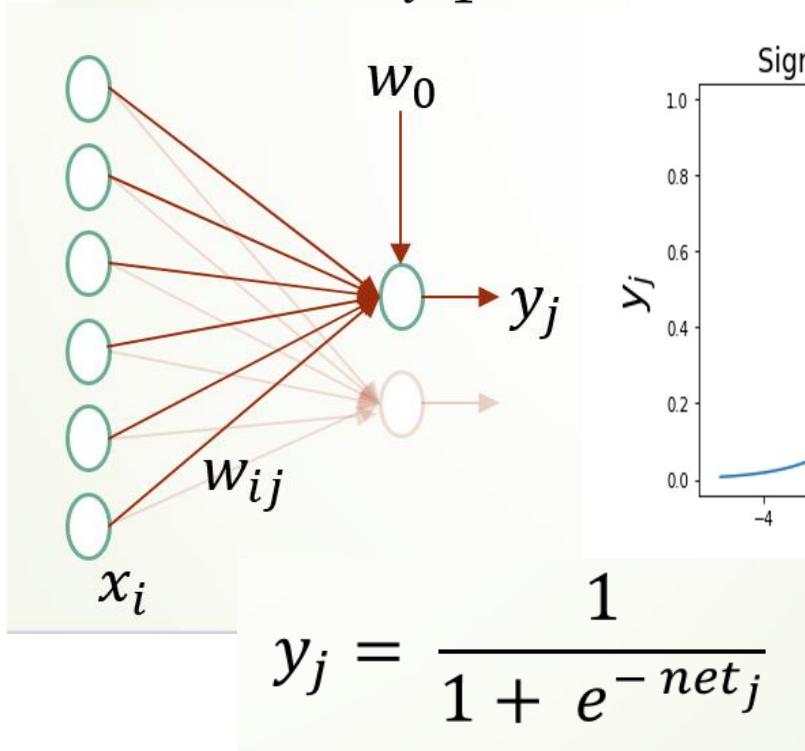
# Biological Inspirations

- Neurons communicate:
  - *electrically* via **action potentials**
  - *chemically* via **neurotransmitters**



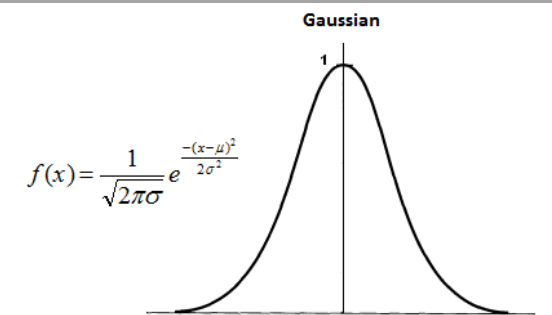
# Artificial Neural Networks

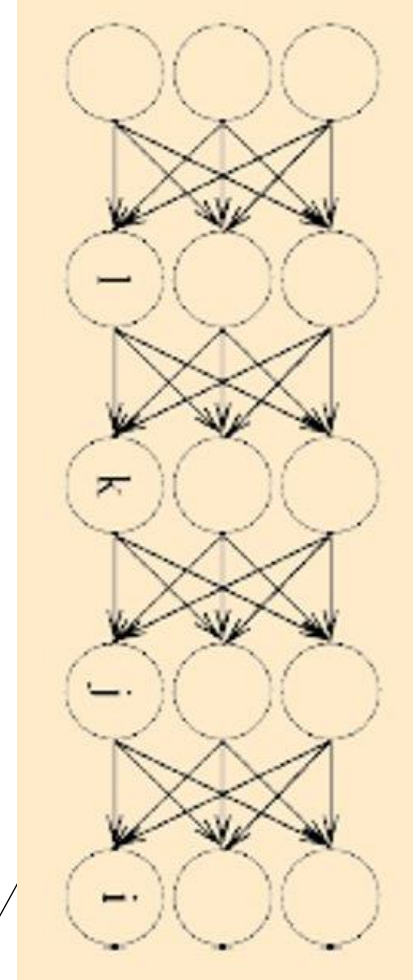
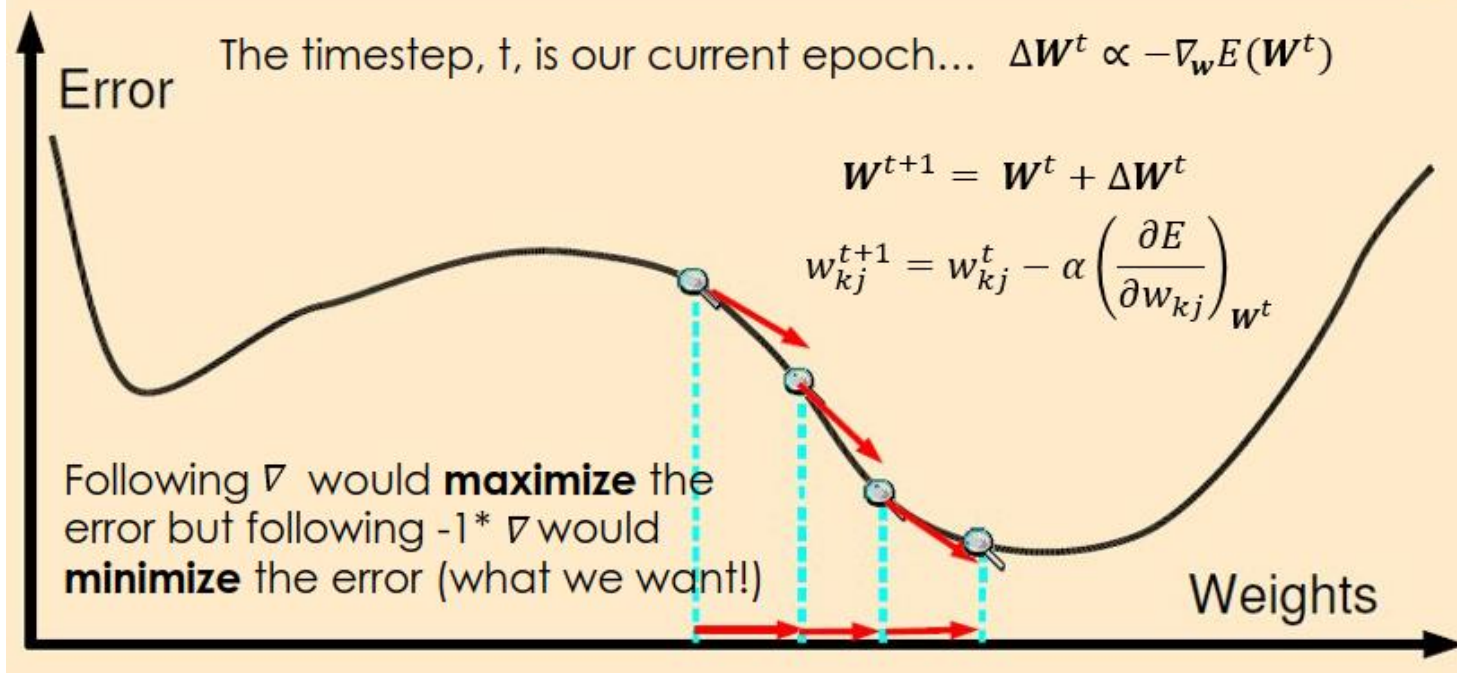
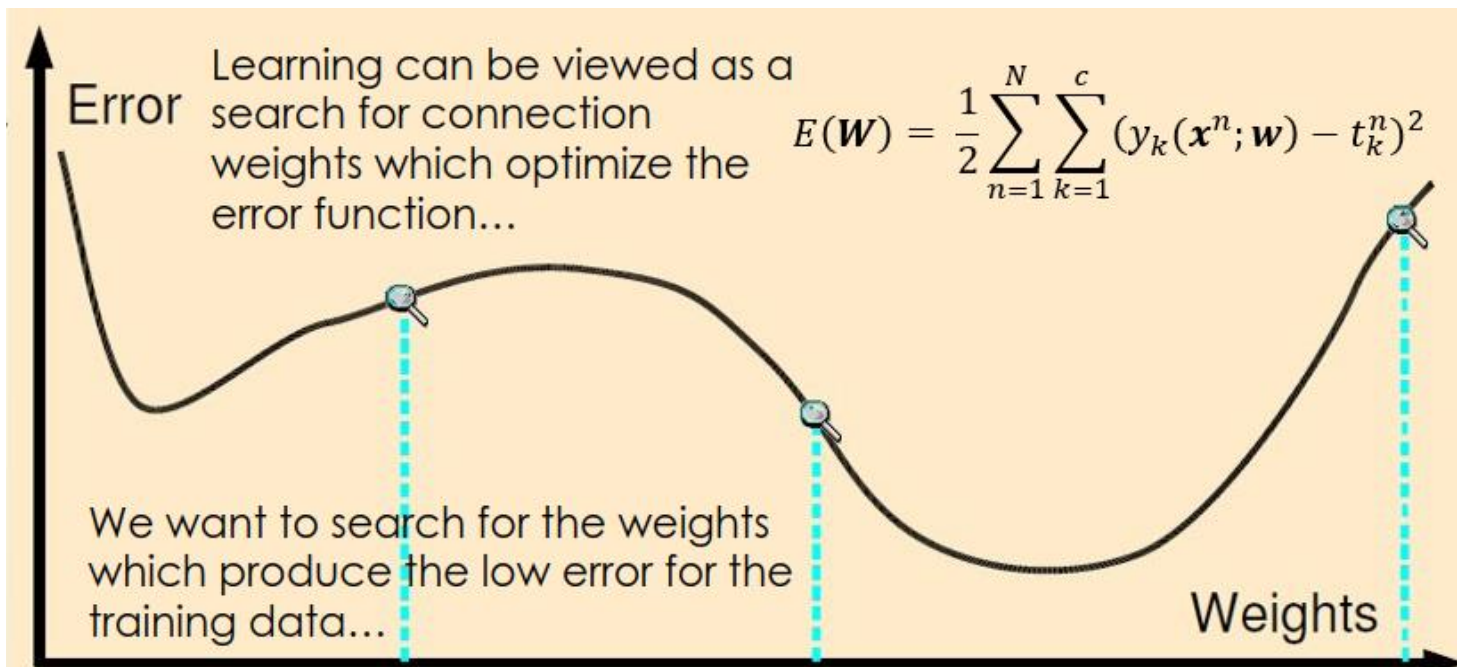
$$net_j = w_0 + \sum_{i=1}^n x_i w_{ij}$$



Feedforward Network

Something for later:  
Radial Basis  
Activation Functions

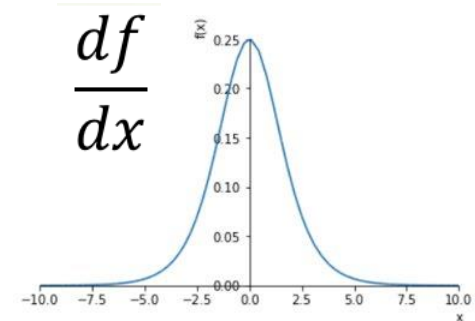




$$\delta_j = g'(net_j) \sum_i w_{ij} \delta_i$$

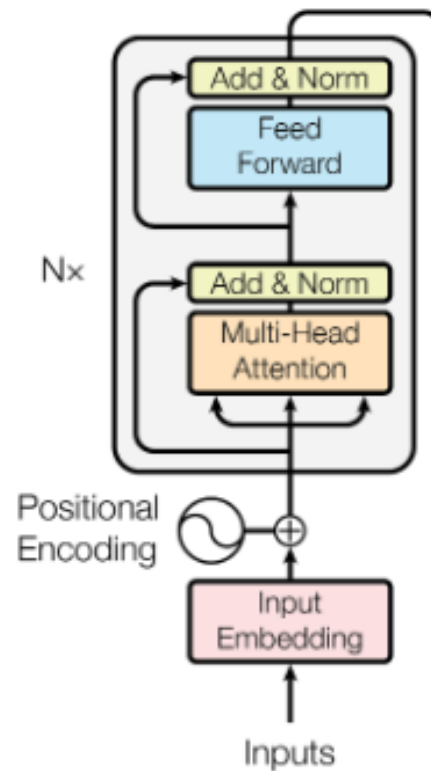
$$\delta_k = g'(net_k) \sum_j w_{jk} \delta_j$$

$$\delta_l = g'(net_l) \sum_k w_{kl} \delta_k$$



# Deep Learning

Vaswani et al., 2017



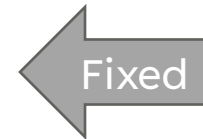
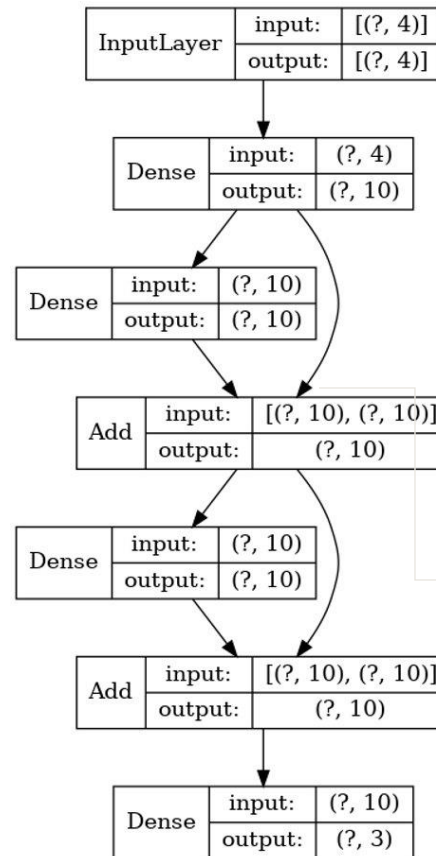
Transformer Architecture

Nonlinear Residual

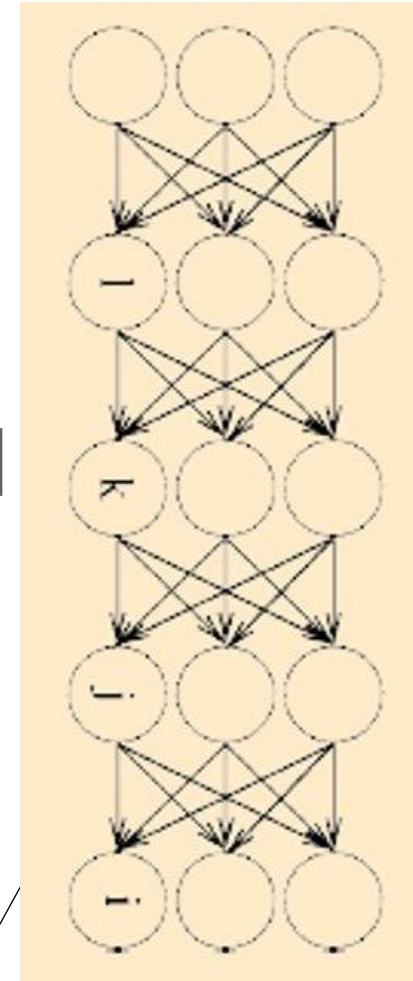
Nonlinear Residual



Number of layers is not as important as **deep residual structure** to say that the model is a form of **deep learning**



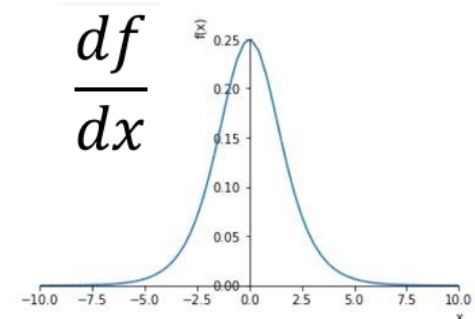
He et al. 2015



$$\delta_j = g'(net_j) \sum_i w_{ij} \delta_i$$

$$\delta_k = g'(net_k) \sum_j w_{jk} \delta_j$$

$$\delta_l = g'(net_l) \sum_k w_{kl} \delta_k$$



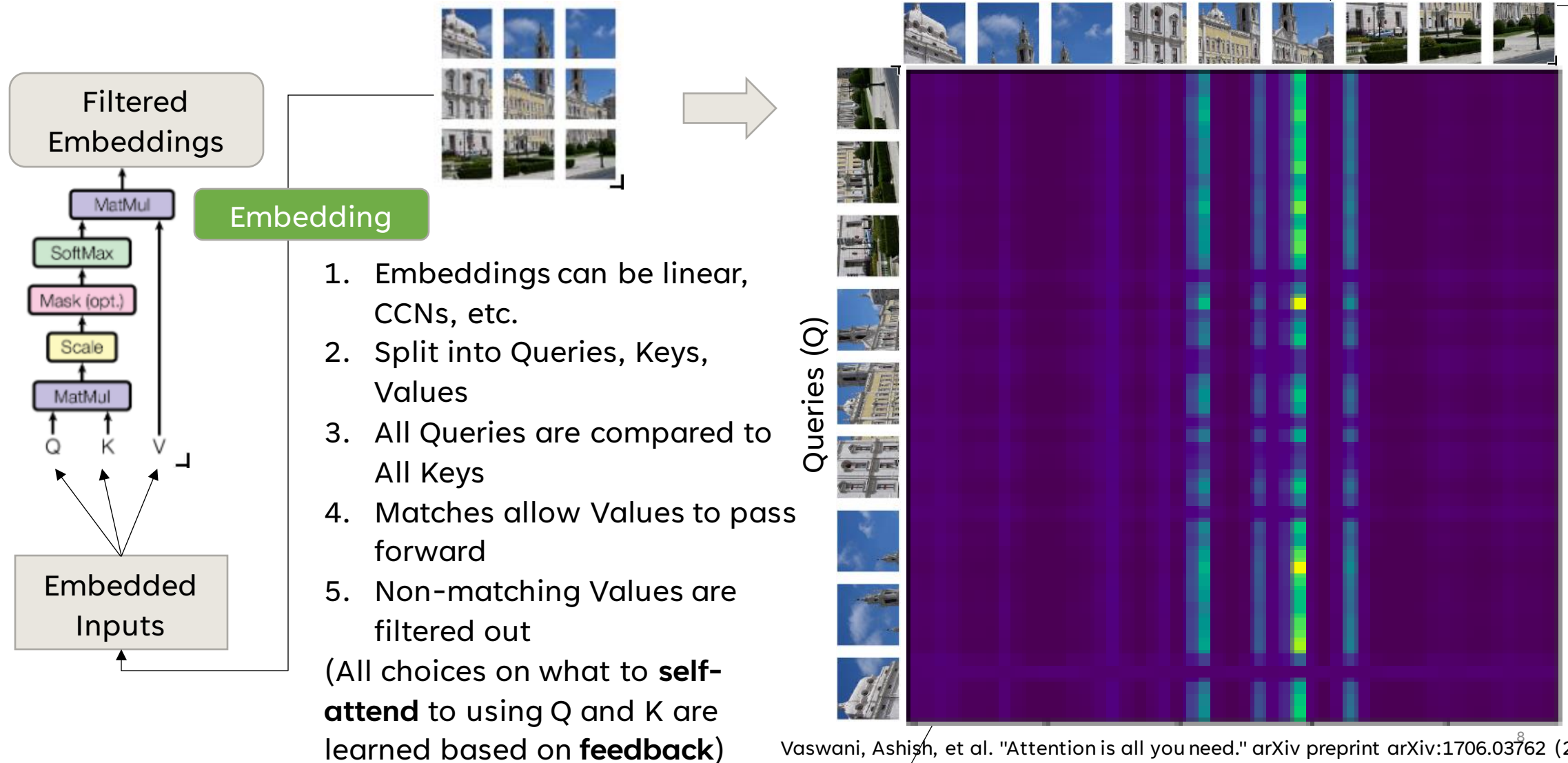


# TRANSFORMER: WHY ATTENTION HELPS



Dosovitskiy, Alexey, et al. "An image is worth 16x16 words: Transformers for image recognition at scale." arXiv preprint arXiv:2010.11929 (2020).

# TRANSFORMER: SELF-ATTENTION MECHANISM



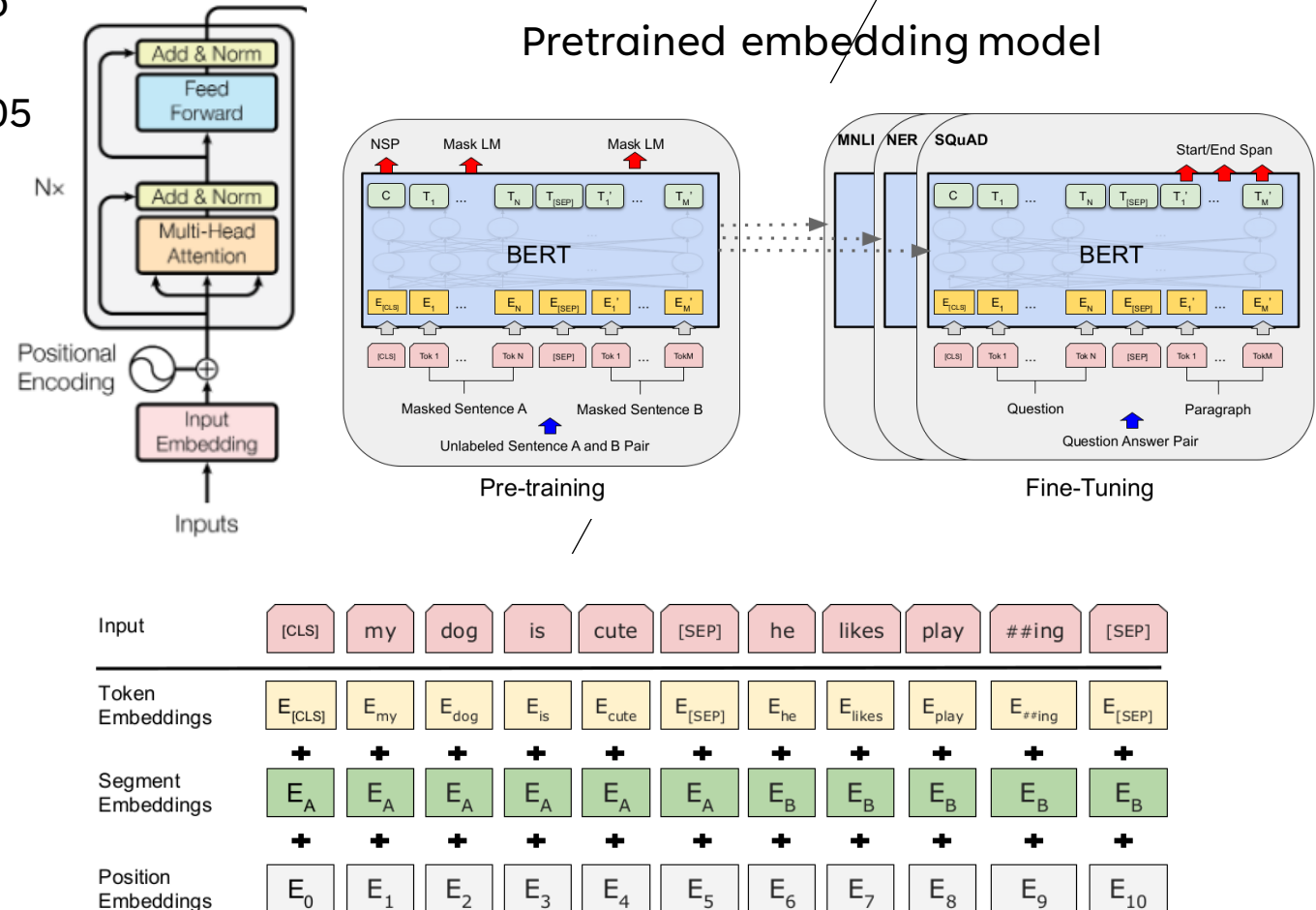
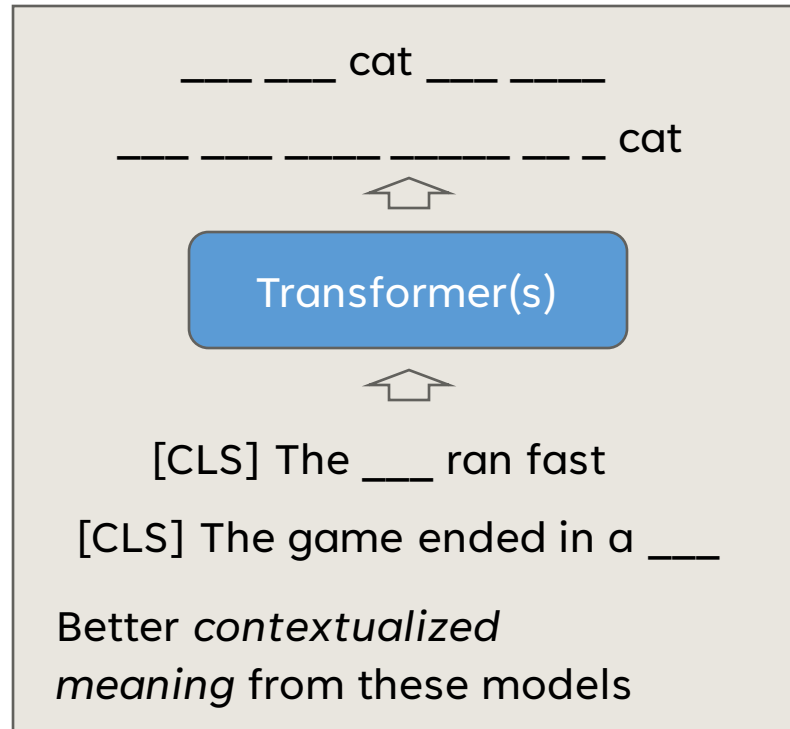


# BERT PRETRAINING METHODOLOGY

Devlin, Jacob, et al. "Bert: Pre-training of deep bidirectional transformers for language understanding." arXiv preprint arXiv:1810.04805 (2018).

NLP Examples:

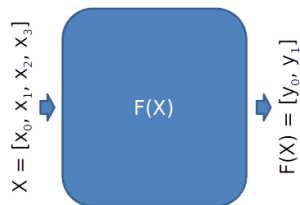
**BERT, mBERT, RoBERTa**



# GPT (GENERATIVE PRETRAINED TRANSFORMER)

- Vaswani et al., 2017 - **Transformer** architecture
- Radford et al., 2018 and Brown et al., 2020
- Simple *generative training* and *testing* procedure, perfectly suited for the *transformer* architecture.
- Very large model, very large data set

1. The  
2. Cat  
3. Ran  
4. Fast

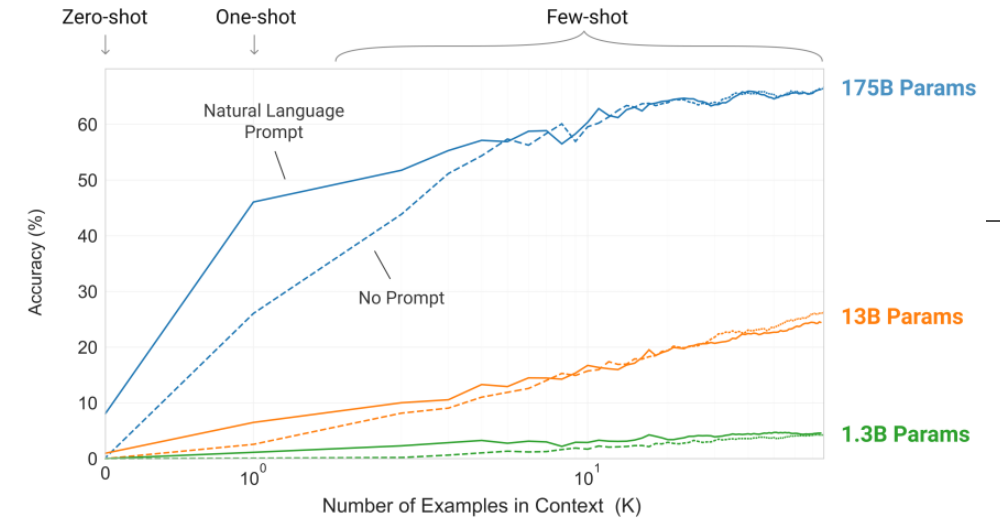


1. Cat  
2. Ran  
3. Fast  
4. <STOP>

The  $[P(\text{duck}), P(\text{cat}), P(\text{fast}), P(\text{no}), \dots]$

The cat  $[P(\text{duck}), P(\text{cat}), P(\text{ran}), \dots]$

The cat ran  $[P(\text{fast}), P(\text{quickly}), P(\text{slowly}), P(\text{no}) \dots]$



GPT-3 (Brown et al. 2020)

[To do a "farduddle" means to jump up and down really fast. An example of a sentence that uses the word farduddle is:]

**One day when I was playing tag with my little sister, she got really excited and she started doing these crazy farduddles.**

[A "yalubalu" is a type of vegetable that looks like a big pumpkin. An example of a sentence that uses the word yalubalu is:]

**I was on a trip to Africa and I tried this yalubalu vegetable that was grown in a garden there. It was delicious.**

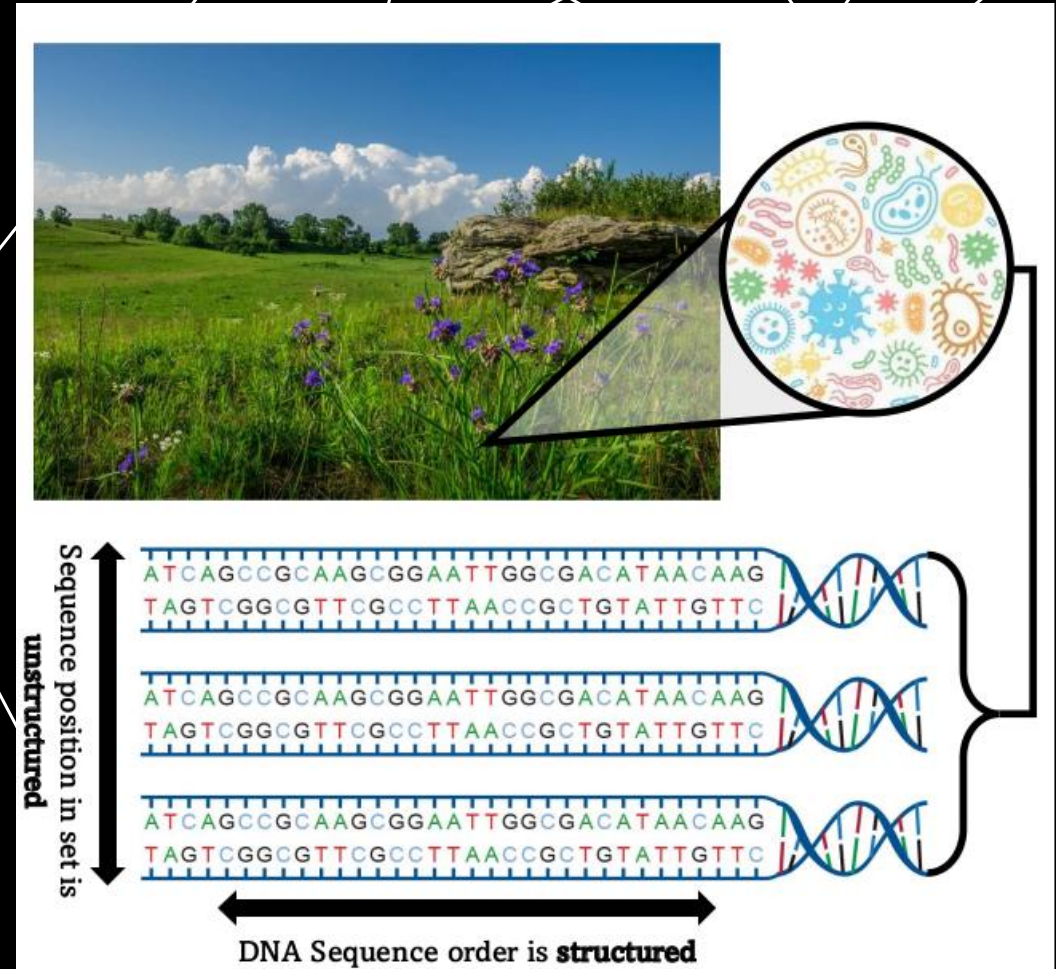
# SETBERT: DEEP LEARNING FOR OMICS DOMAINS

- We began studying high-throughput sequencing data (HTS) in 2020
- NSF Award# 1933925 (**Collaborative Research: RoL: Impacts of plants and communities on soil microbial composition and function across phylogenetic scales**)
- Structured – DNA sequences are *ordered* sequences
- Unstructured – HTS FASTQ are *unordered* sets of sequences



David Ludwig  
COMS Ph.D. Student

```
TACGTAGGGTGCAAGC
AACGTAGGTACCGAGC
TACGAAGGGGCTAGC
TACAGAGGGTGCAAGC
TACGTAGGTGGCAAGC
GACGTAGGGTGCGAGC
...
TACGGAGGGGCTAGC
GACGAACCGTCCGAAC
CACGGACCGCACGAAC
```





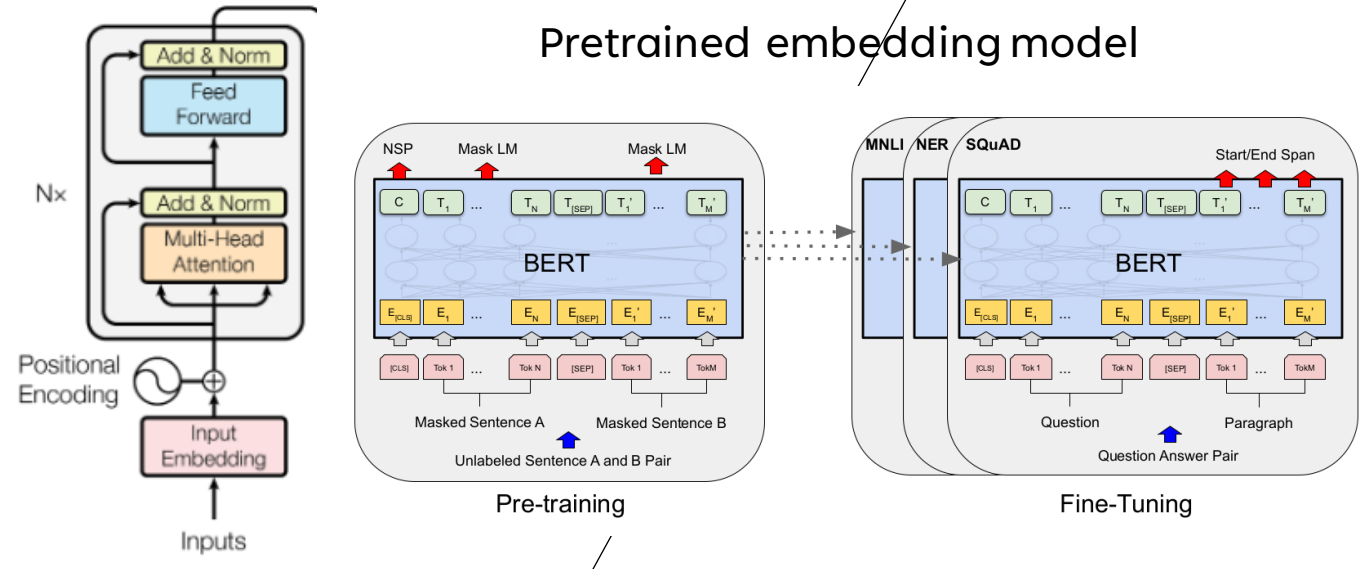
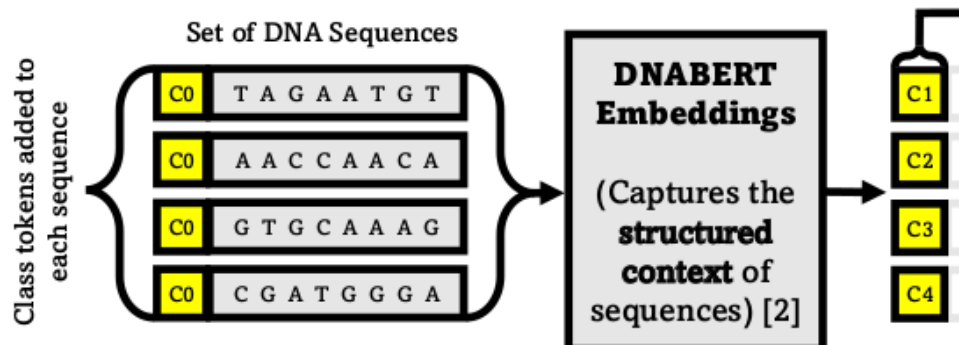
# BERT PRETRAINING METHODOLOGY

Devlin, Jacob, et al. "Bert: Pre-training of deep bidirectional transformers for language understanding." arXiv preprint arXiv:1810.04805 (2018).

NLP Examples:

**BERT, mBERT, RoBERTa**

Analogously applied to HTS data...

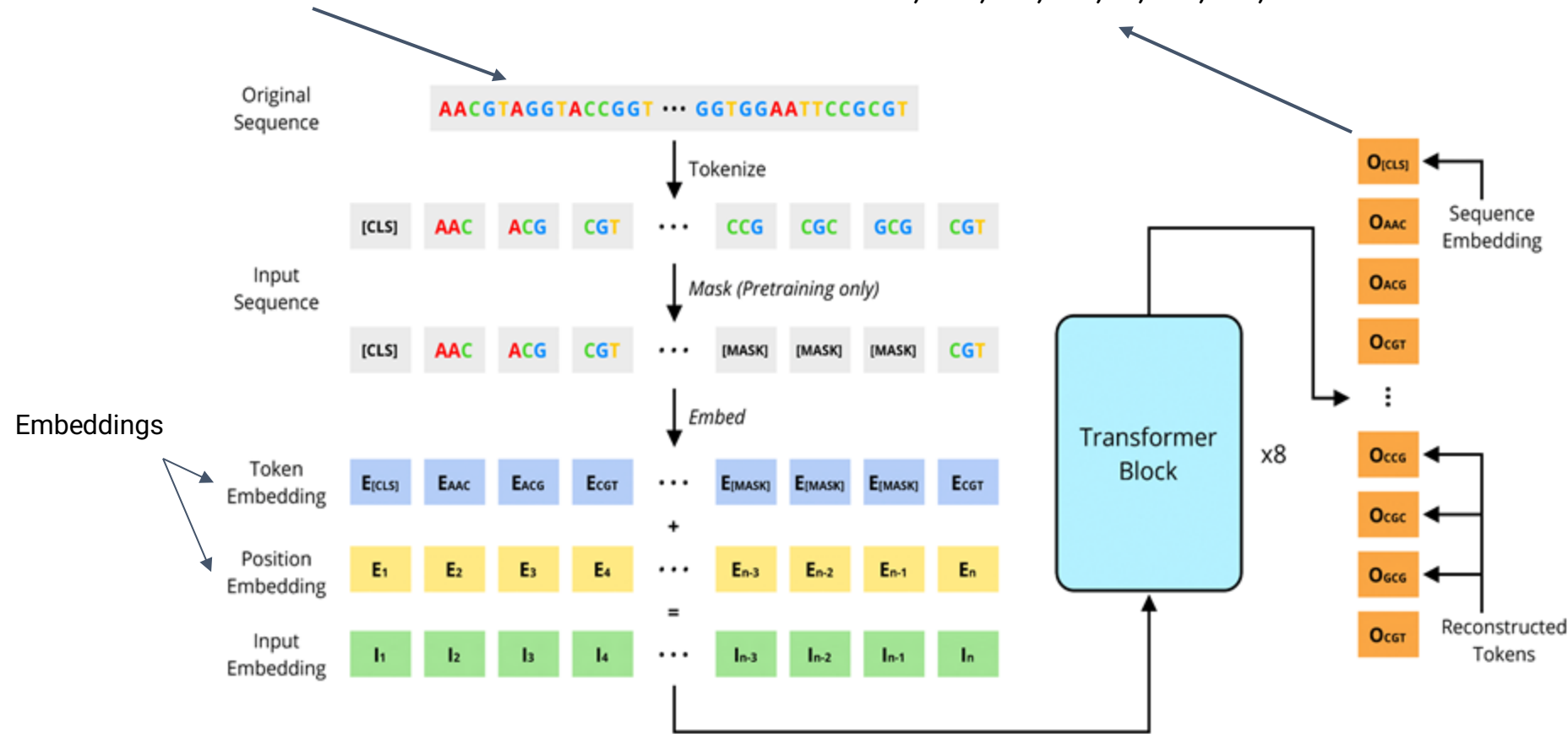


Input	[CLS]	my	dog	is	cute	[SEP]	he	likes	play	##ing	[SEP]
Token Embeddings	$E_{[CLS]}$	$E_{my}$	$E_{dog}$	$E_{is}$	$E_{cute}$	$E_{[SEP]}$	$E_{he}$	$E_{likes}$	$E_{play}$	$E_{\#ing}$	$E_{[SEP]}$
Segment Embeddings	+	+	+	+	+	+	+	+	+	+	+
	$E_A$	$E_A$	$E_A$	$E_A$	$E_A$	$E_A$	$E_B$	$E_B$	$E_B$	$E_B$	$E_B$
Position Embeddings	+	+	+	+	+	+	+	+	+	+	+
	$E_0$	$E_1$	$E_2$	$E_3$	$E_4$	$E_5$	$E_6$	$E_7$	$E_8$	$E_9$	$E_{10}$

# DNABERT: SINGLE SEQUENCE FOUNDATIONS

Map raw DNA sequence to a high-dimensional vector embedding

AACGTAGGTACCGGT...GGTGGAATTCCGCGT  $\longrightarrow$   $\langle 0.1, -0.3, 0.2, 0.2, \dots, 0.4, 0.9, -0.2 \rangle$

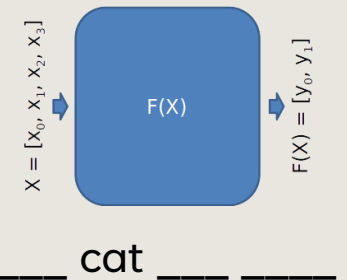


## BERT Training Methodology

- [Vaswani et al., 2017](#) - Transformer architecture
- [Devlin et al., 2018](#)

The \_\_\_\_ ran fast

The game ended in a \_\_\_\_

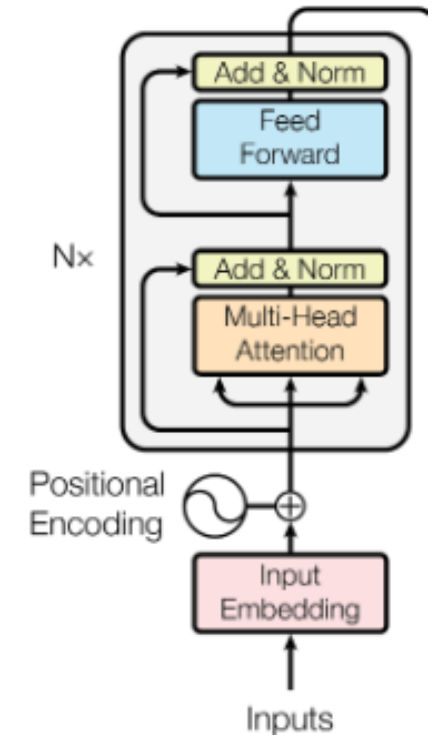
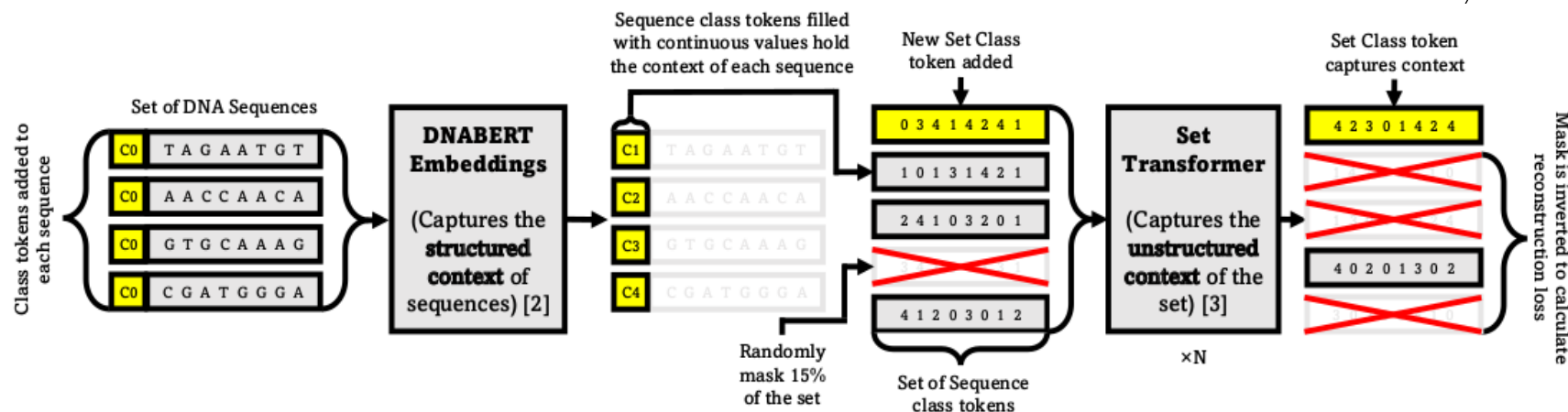


cat  
Better contextualized meaning from these models

Ji et al., *Bioinformatics*, Volume 37, Issue 15, August 2021, Pages 2112–2120, <https://doi.org/10.1093/bioinformatics/btab083>

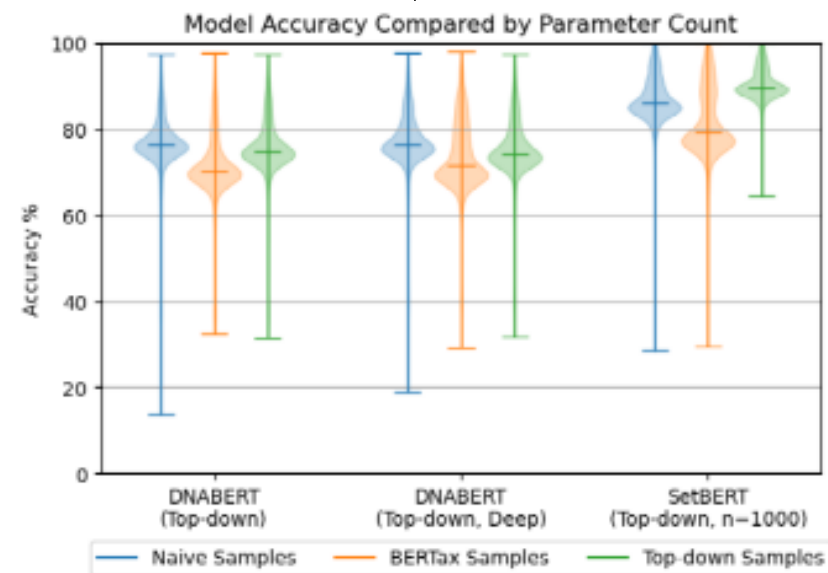
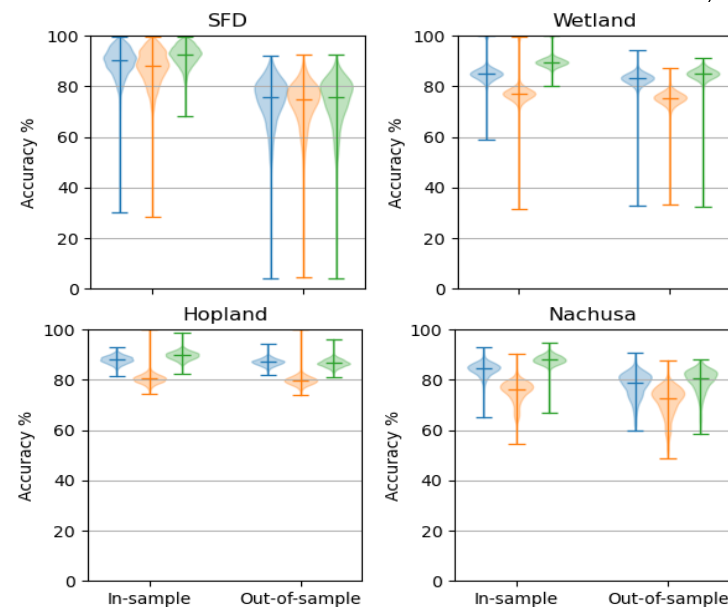
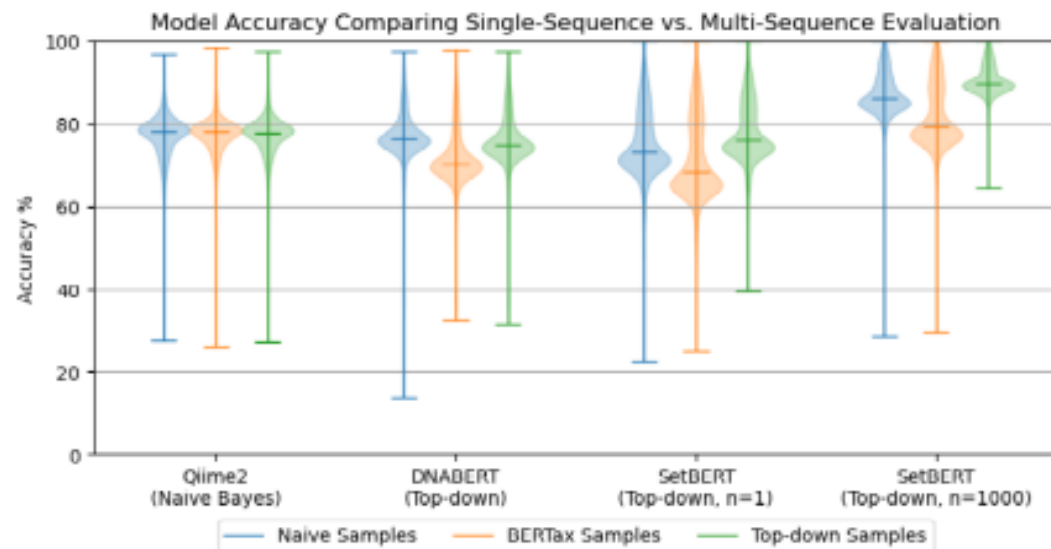
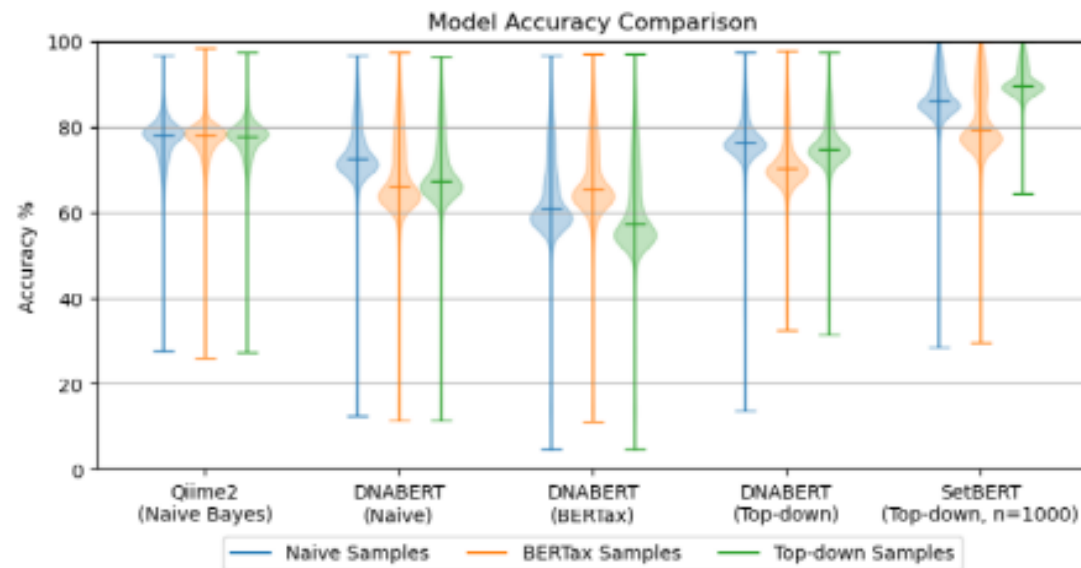
# SETBERT: DNABERT + SET TRANSFORMER (LEE ET AL., 2019)

- The Set Transformer component integrates single-sequence embeddings in a *permutation-equivariant* manner
- Removal of *position encodings* is all that is needed to make the MHA module of Transformers perform in this manner
- We construct an HTS sample as a set of embeddings provided by our pretrained DNABERT module
- Pretraining can be performed for sets of sequences using fixed DNABERT embeddings
- Fine-tuning can then be performed for any downstream task
- Changes in attention between pretrained embeddings (or compared to fine-tuned results) can be made, and attention attribution is also compatible

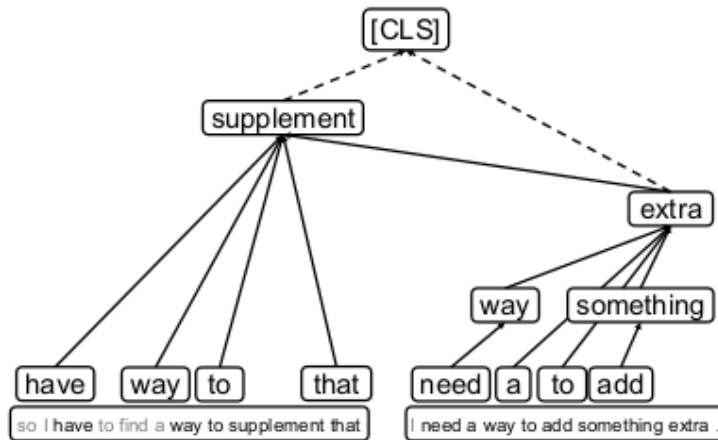




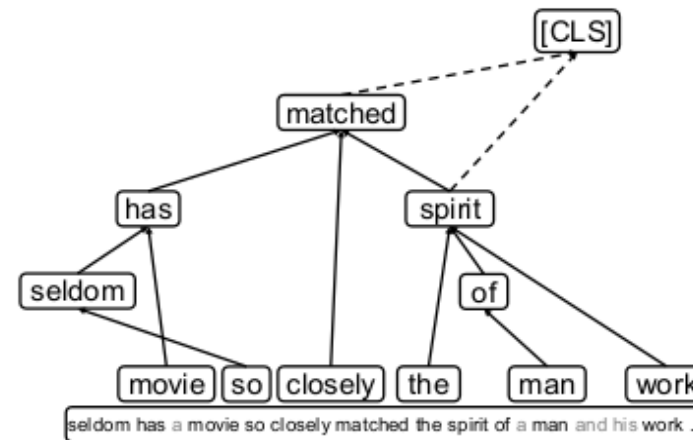
# TAXONOMIC PREDICTION: GENUS LEVEL



# INTERPRETING ATTENTION MECHANISMS



(a) Example from MNLI



(b) Example from SST-2

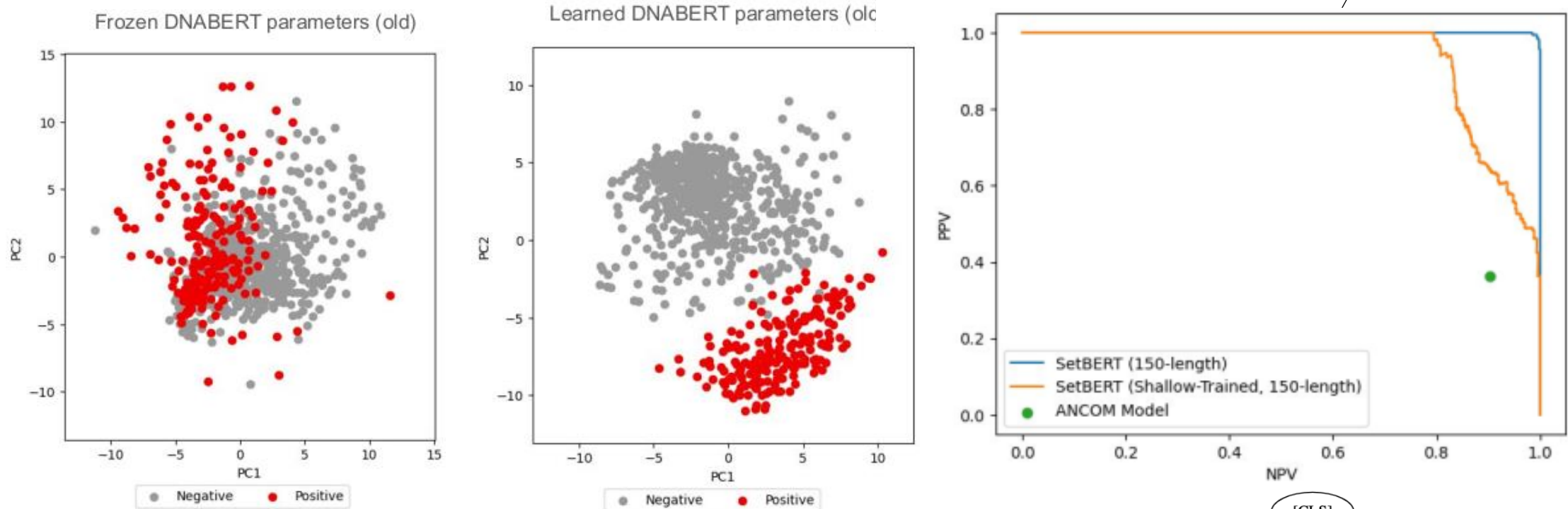
Figure 5: Examples of attribution trees. (a) is from MNLI, which is predicted as entailment by BERT. (b) is from SST-2, which is predicted as positive by BERT. The grey words from the inputs do not appear in the attribution trees.



(b) Attribution Score

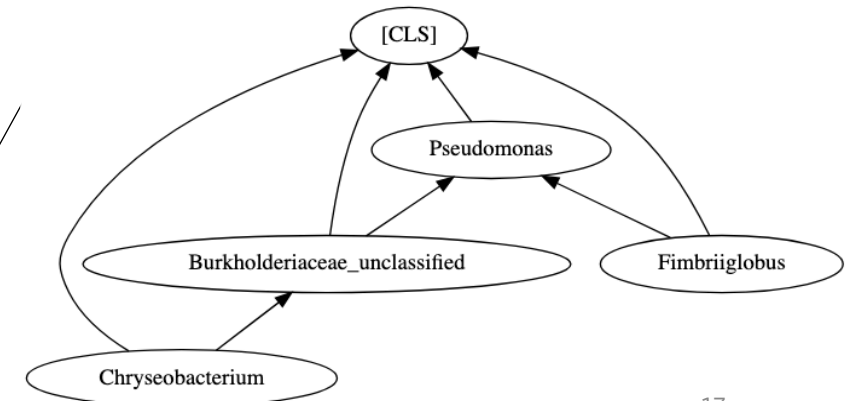
Hao, Yaru, et al. "Self-attention attribution: interpreting information interactions inside transformer." arXiv preprint arXiv:2004.11207 (2020).

# SNAKE FUNGAL DISEASE (PATHOGEN PRESENCE/ABSENCE)



Improvements due to allowing gradient flow back into the DNABERT component of the SetBERT models. (Previous model used fixed DNABERT embeddings to preserve GPU memory).

Collaboration with Walker Lab at MTSU (Donald Walker, Department of Biology, MTSU)

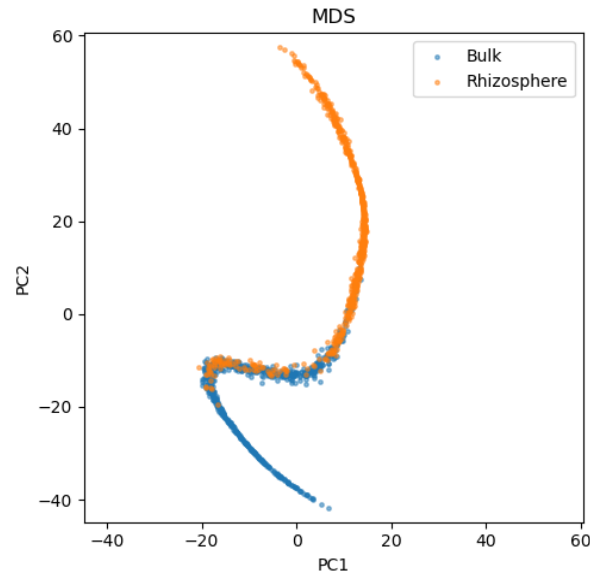
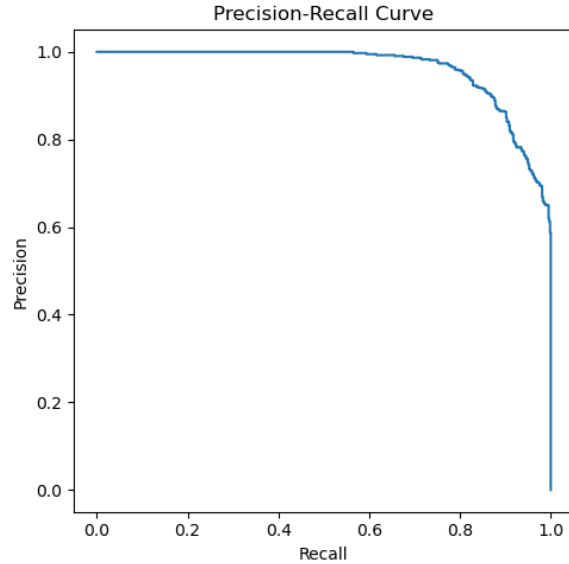




# ATTENTION SHIFTS (PRETRAINED TO FINE-TUNED)

Predicted Pathogen Presence	Taxa Category	Phylum	Class	Order	Family	Genus	Score
Positive	Highest Scoring	Actinobacteria	Actinobacteria	Corynebacteriales	Dietziaceae	Dietzia	25,422
		Bacteroidetes	Bacteroidia	Flavobacteriales	Weeksellaceae	Chryseobacterium	25,348
		Proteobacteria	Alphaproteobacteria	Sphingomonadales	Sphingomonadaceae	Sphingomonas	22,800
		Firmicutes	Bacilli	Bacillales	Staphylococcaceae	Staphylococcus	22,000
		Proteobacteria	Betaproteobacteria	Burkholderiales	Comamonadaceae	Comamonas	21,500
		Proteobacteria	Gammaproteobacteria	Enterobacteriales	Enterobacteriaceae	Escherichia	21,001
		Firmicutes	Bacilli	Lactobacillales	Lactobacillaceae	Lactobacillus	20,500
		Bacteroidetes	Bacteroidia	Bacteroidales	Porphyromonadaceae	Parabacteroides	20,250
		Actinobacteria	Actinobacteria	Actinomycetales	Actinomycetaceae	Actinomyces	20,001
		Firmicutes	Clostridia	Clostridiales	Ruminococcaceae	Ruminococcus	19,850
	Lowest Scoring	Actinobacteria	Actinobacteria	Micrococcales	Dermacoccaceae	Flexivirga	-7,073
		Proteobacteria	Gammaproteobacteria	Enterobacteriales	Enterobacteriaceae	Enterobacteriaceae_unclassified	-6,149
		Bacteroidetes	Bacteroidia	Cytophagales	Microscillaceae	Siphonobacter	-4,966
		Proteobacteria	Gammaproteobacteria	Pseudomonadales	Moraxellaceae	Acinetobacter	-3,918
		Firmicutes	Bacilli	Lactobacillales	Enterococcaceae	Enterococcus	-3,339
		Actinobacteria	Actinobacteria	Micrococcales	Dermabacteraceae	Brachybacterium	-2,196
		Proteobacteria	Gammaproteobacteria	Pseudomonadales	Pseudomonadaceae	Pseudomonas	-2,130
		Firmicutes	Bacilli	Bacillales	Bacillales_unclassified	Bacillales_unclassified	-1,825
		Firmicutes	Erysipelotrichia	Erysipelotrichales	Erysipelotrichaceae	Erysipelothrix	-1,765
		Firmicutes	Bacilli	Bacillales	Staphylococcaceae	Staphylococcus	-1,346

# RILL K SOIL VS. RHIZOSPHERE DISTINCTIONS



Actinobacteriota  
Acidobacteriota  
Planctomycetota  
Actinobacteriota  
Gemmatimonadota  
Proteobacteria  
Acidobacteriota  
Proteobacteria  
Actinobacteriota  
Acidobacteriota

Firmicutes  
Firmicutes  
Firmicutes  
Firmicutes  
Proteobacteria  
Bacteroidota  
Verrucomicrobiota  
Firmicutes  
Bacteroidota  
Bacteroidota

Thermoleophilia  
Vicinamibacteria  
Planctomycetes  
Thermoleophilia  
Gemmatimonadetes  
Alphaproteobacteria  
Thermoanaerobaculia  
Alphaproteobacteria  
Thermoleophilia  
Vicinamibacteria

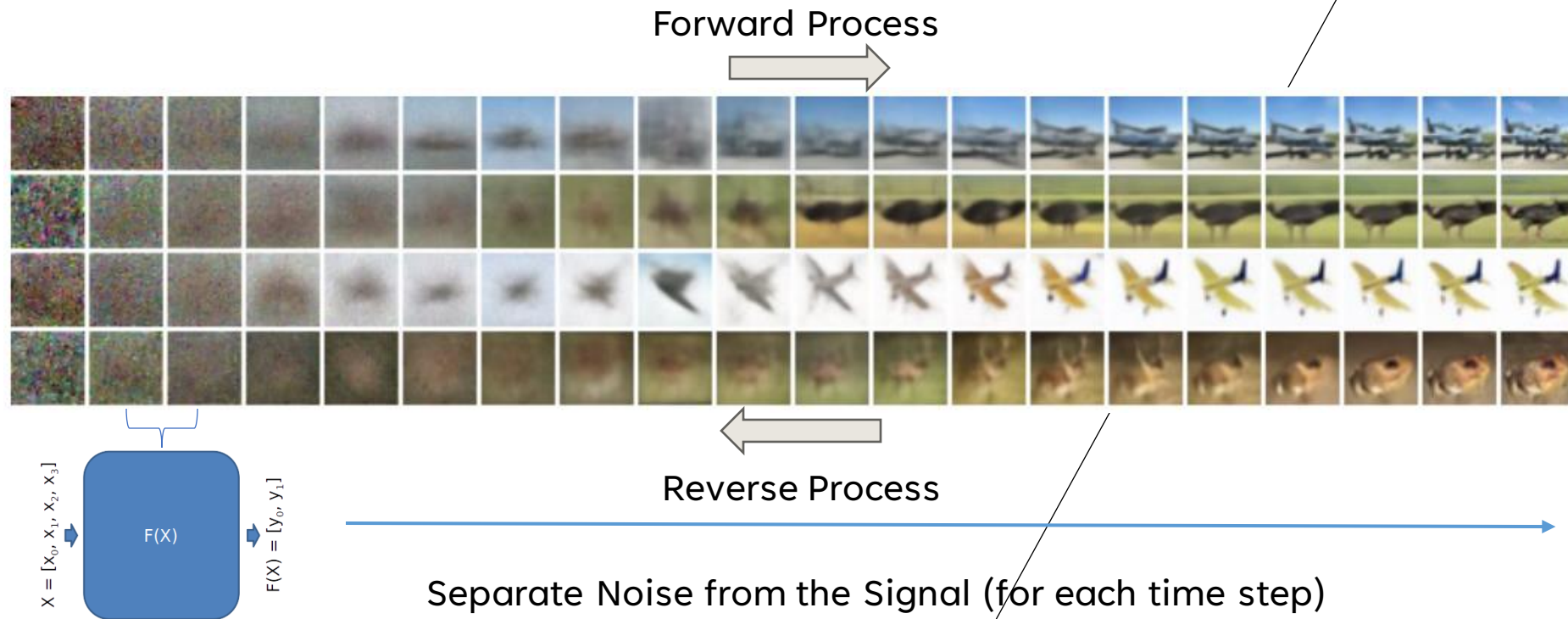
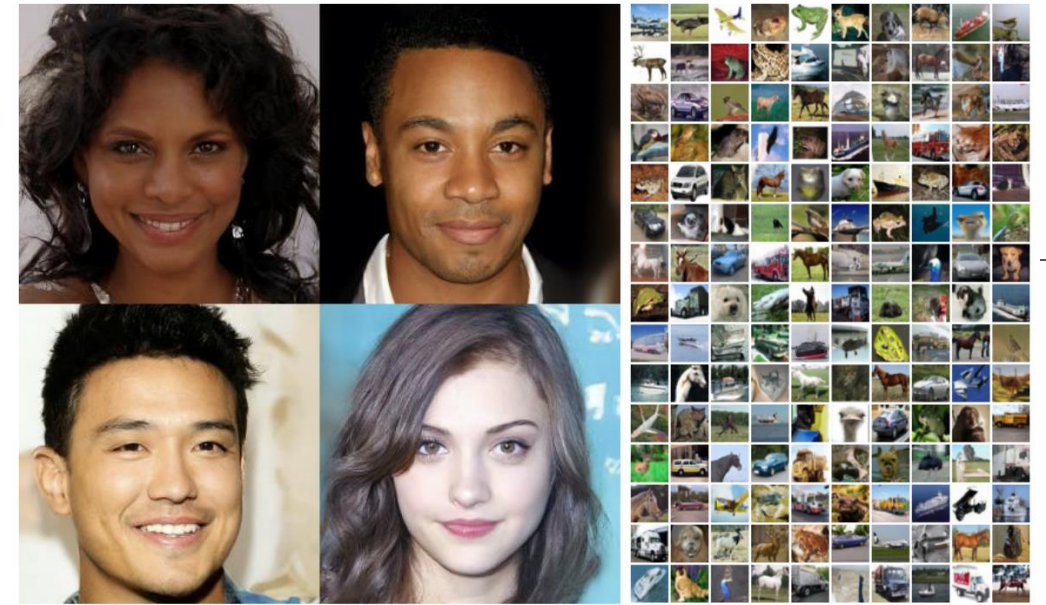
Bacilli  
Bacilli  
Bacilli  
Bacilli  
Alphaproteobacteria  
Bacteroidia  
Verrucomicrobiae  
Bacilli  
Bacteroidia  
Bacteroidia

Solirubrobacterales  
Vicinamibacterales  
Gemmatales  
Solirubrobacterales  
Gemmatimonadales  
Rhizobiales  
Thermoanaerobaculales  
Rhizobiales  
Gaiellales  
Vicinamibacterales

Bacillales  
Bacillales  
Bacillales  
Bacillales  
Sphingomonadales  
Chitinophagales  
Chthoniobacterales  
Bacillales  
Chitinophagales  
Chitinophagales

# DIFFUSION MODELS

- Ho et al., 2020
- Denoising Diffusion Probabilistic Models (DDPM)
- An elegant solution to the *mode collapse* issue with generative modeling tasks

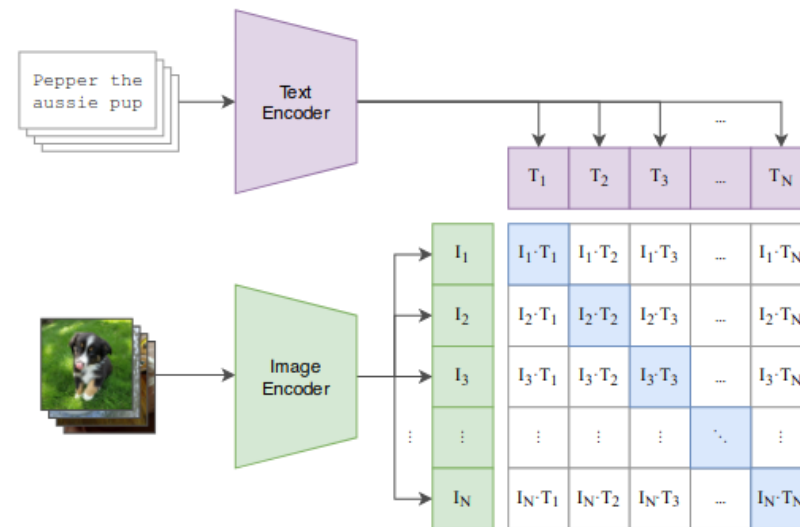




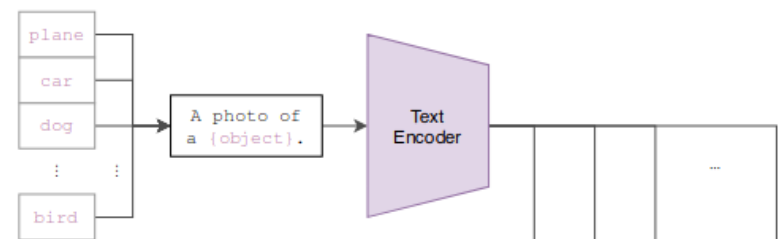
# CONTRASTIVE EMBEDDING ALIGNMENT ACROSS DOMAINS

- [Radford et al., 2021](#)
- Contrastive learning can be done to align the learned information from one domain to another related domain
- Image to text
- Text to image
- Doesn't have to just be these domains...
- Zero-shot prediction still possible, just like with the earlier SimCLR approach

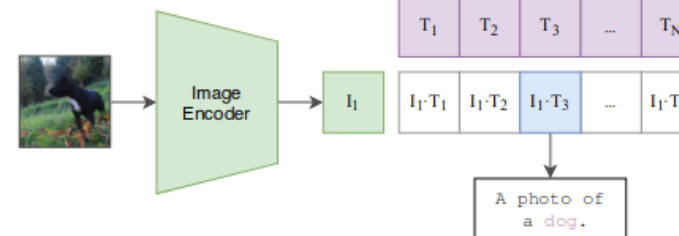
## (1) Contrastive pre-training



## (2) Create dataset classifier from label text



## (3) Use for zero-shot prediction



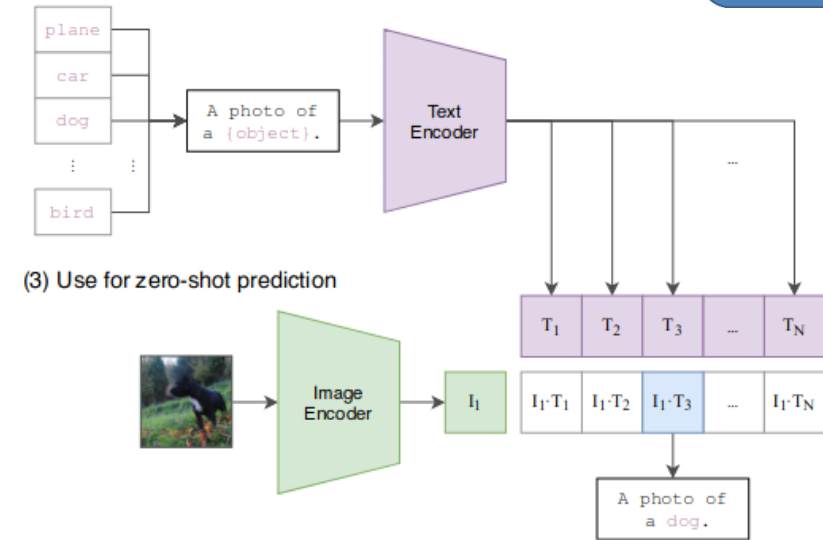
# CONDITIONAL GENERATION WITH CONTRASTIVE EMBEDDINGS: DALL-E

- Ramesh et al., 2022
- Reverse diffusion process can be trained while including a *contrastive embedding*
- The text encoder can generate a *similar* contrastive embedding
- The diffusion model can use the text's contrastive embedding to extract a *similar* image
- Natural variation in the DDPM model and the large variation in training data allows for creative, generative modeling

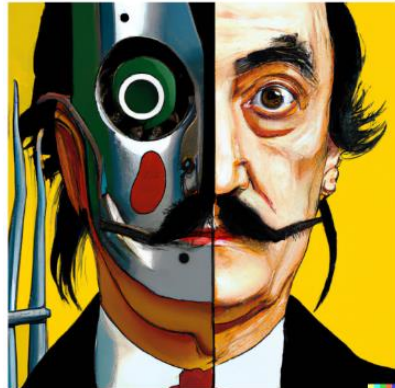
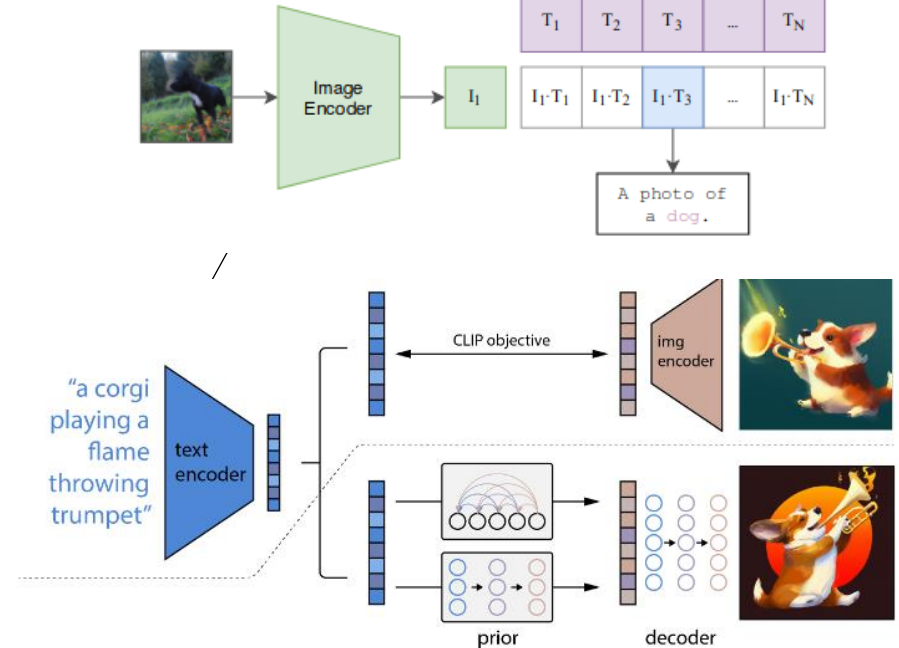


## Conditioned Diffusion

(2) Create dataset classifier from label text



(3) Use for zero-shot prediction



vibrant portrait painting of Salvador Dalí with a robotic half face



a shiba inu wearing a beret and black turtleneck



a close up of a handpalm with leaves growing from it



an espresso machine that makes coffee from human souls, artstation



panda mad scientist mixing sparkling chemicals, artstation



a corgi's head depicted as an explosion of a nebula

# SOME FINAL THOUGHTS AND TAKE-AWAYS

1. Language, images, audio (not covered here, but improvements in this domain are steady now as well): *all human communication* and therefore all domains of human knowledge are directly impacted by Generative AI.
2. Recent work has shown they often out-perform humans on abstract reasoning tasks ([Webb et al., Aug. 2023](#)). "*Our results indicate that large language models such as GPT-3 have acquired an emergent ability to find zero-shot solutions to a broad range of analogy problems.*"
3. Entire scientific process pipelines are being assembled using multiple coordinating generative models: hypothesis generation, experimental design, search and execution, and analysis/revision.
4. What about Quantum?
  1. Some functions of interest rely on quantum phenomena and AI may soon be used to learn these quantum functions instead of relying on human programming.
  2. Better optimization and approximation may be possible by integrating quantum algorithms into AI systems and this is an active area of research.

A series of white, overlapping geometric lines and polygons on a black background, located on the left side of the slide.

# THANK YOU

Joshua L. Phillips

[Joshua.Phillips@mtsu.edu](mailto:Joshua.Phillips@mtsu.edu)

<https://phillips-lab.org/>