# CS 839: Foundation Models
# **Course Overview**

Fred Sala

University of Wisconsin-Madison

**Sept. 4, 2025**

# **Logistics**: Lecture Location

- In-person in **Morgridge 2532**
  - Will have slides / blackboard usage
  - Blackboard for theory; slides for model diagrams etc.

- Planning to record---final decision TBD.

# **Logistics**: Enrollment

- Currently at capacity, approx. 75 students

  - Some folks on waitlist may not make it in
  - Decent chance many of the waitlist folks will

- **Sorry** ☹ … will be offered again

# **Logistics**: Teaching Team

Instructor: **Fred Sala**
- Location: 5514 Morgridge Hall
- Office Hours: TBD

TA: **Changho Shin**
- Location: 5548 Morgridge Hall
- Office Hours: TBD

- Note: times possibly **subject to change**

# **Logistics**: Content

Three locations:

- 1. **Course website**: https://pages.cs.wisc.edu/~fredsala/cs839/fall2025/

- 2. **Piazza**.  https://piazza.com/wisc/fall2025/d89e
  - access code: *introtofm*
  - **Preferred for questions!**

- 3. **Canvas**



Flickr Photo by Jacob Surland - creative commons licensed ( BY-NC )

# Course Content / Schedule

| Date | Topic | Readings | Other |
|---|---|---|---|
| Thursday Sept. 4 | Introduction and Course Overview | | |
| Tuesday Sept. 9 | Machine Learning Mini-Review | • Patterns, Predictions, and Actions | |
| Thursday Sept. 11 | Architectures I: Transformers & Attention | • Attention Is All You Need<br>• The Illustrated Transformer<br>• Tokenization<br>• RoFormer: Rotary Position Embedding (RoPE)<br>• ALiBi: Train Short, Test Long | |
| Tuesday Sept. 16 | Architectures II: Subquadratic Architectures | • Mamba: Linear-Time Sequence Modeling with Selective State Spaces<br>• Mamba: The Hard Way<br>• Retentive Network (RetNet)<br>• Transformers are SSMs | |
| Thursday Sept. 18 | Models I: Encoder-Decoder, Encoder-Only | • BERT Paper<br>• RoBERTa Paper<br>• T5 Paper | HW 1 Release |
| Tuesday Sept. 23 | Models II: Decoder-Only. Start Prompting I: Basics | • GPT-3 Paper<br>• Llama 3.1 Paper | |
| Thursday Sept. 25 | Prompting I: Basics, In-Context Learning, Chain-of-Thought | • Finetuned Language Models Are Zero-Shot Learners<br>• Learning without training: The implicit dynamics of in-context learning<br>• CoT<br>• Large Language Models are Zero-Shot Reasoners | |
| Tuesday Sept. 30 | Specialization: Fine-Tuning, Adaptation, Editing | • Low-Rank Adaptation of Large Language Models<br>• Fast Model Editing at Scale<br>• QLoRA: Efficient Finetuning of Quantized LLMs | |
| Thursday Oct. 2 | Alignment: RHLF, DPO, and friends | • Hugging Face RLHF Blog<br>• Training language models to follow instructions with human feedback<br>• Direct Preference Optimization: Your Language Model is Secretly a Reward Model<br>• Constitutional AI: Harmlessness from AI Feedback<br>• RLAIF: Scaling RLHF with AI Feedback | Presentation I Release |
| Tuesday Oct. 7 | Reinforcement Learning with Verifiable Rewards (RLVR) | • RLVR (Chapter 3.2) | HW 2 Release |
| Thursday Oct. 9 | Efficient Training | • FlashAttention<br>• Megatron-LM | |
| Tuesday Oct. 14 | Efficient Inference | • Fast Inference from Transformers via Speculative Decoding<br>• Efficient Memory Management for Large Language Model Serving with PagedAttention<br>• LLM.int8()<br>• AWQ Quantization<br>• Medusa Decoding | |
| Thursday Oct. 16 | Evaluation, Benchmarks, and LLM-as-a-Judge, | • HELM: Holistic Evaluation of Language Models<br>• MMLU-Pro<br>• LLM-as-a-Judge: MT-Bench & Chatbot Arena | Project Information Out |
| Thursday Oct. 23 | Reasoning | • Self-Consistency Improves Chain-of-Thought<br>• Tree of Thoughts<br>• Let's Verify Step by Step (Process Supervision)<br>• Program-of-Thoughts Prompting<br>• OpenAI o1 System Card (Reasoning models)<br>• DeepSeek-R1: Incentivizing Reasoning via RL | |
| Tuesday Oct. 21 | Agents | • ReAct: Reasoning + Acting<br>• Toolformer<br>• Reflexion: Verbal RL for Agents<br>• AgentBench: Evaluating LLMs as Agents<br>• OSWorld: Multimodal Computer Use Benchmark<br>• WebArena: Realistic Web Environment<br>• Voyager: Embodied Agent in Minecraft | |
| Tuesday Oct. 28 | Multimodal Foundation Models | • CLIP: Learning Transferable Visual Models<br>• Flamingo<br>• LLaVA: Visual Instruction Tuning<br>• SAM 2: Segment Anything in Images & Videos<br>• MMMU: Multimodal Understanding & Reasoning<br>• AlphaFold<br>• ClimaX | HW 3 Release |
| Thursday Oct. 30 | Scaling & Scaling Laws | • "Compute-Optimal" LLMs<br>• Scaling Laws for Neural LMs | HW2 Due |
| Tuesday Nov. 4 | Diffusion Models | • Denoising Diffusion Probabilistic Models<br>• Maximum Likelihood Training of Score-Based Diffusion Models<br>• Diffusion language models<br>• Score-based Generative Modeling via SDEs<br>• Flow Matching for Generative Modeling<br>• Consistency Models | |
| Thursday Nov. 6 | Security, Privacy, Toxicity + Future Areas | • Extracting Training Data from LLMs<br>• LLMs Can Be Strong Differentially Private Learners<br>• Red Teaming Paper | |

Launchpad

# **Logistics**: Lecture Formats

Two types of class sessions:

- **Type 1: Lectures**
  - Mostly slides, some whiteboard
  - Will take some breaks, 1-2 during the lecture
  - Can ask questions---during lecture and breaks

- **Type 2: Paper Presentations**
  - More info on later slides.

- Start with Type 1, conclude semester with Type 2

# **Logistics**: Assignments & Grades

**Homeworks**:
- 3-4 or so, worth 30% total
- Posted after class; due when class starts on due date. About 2-3 weeks given for each one

**Class Presentation**:
- Total of 30%. Present a paper
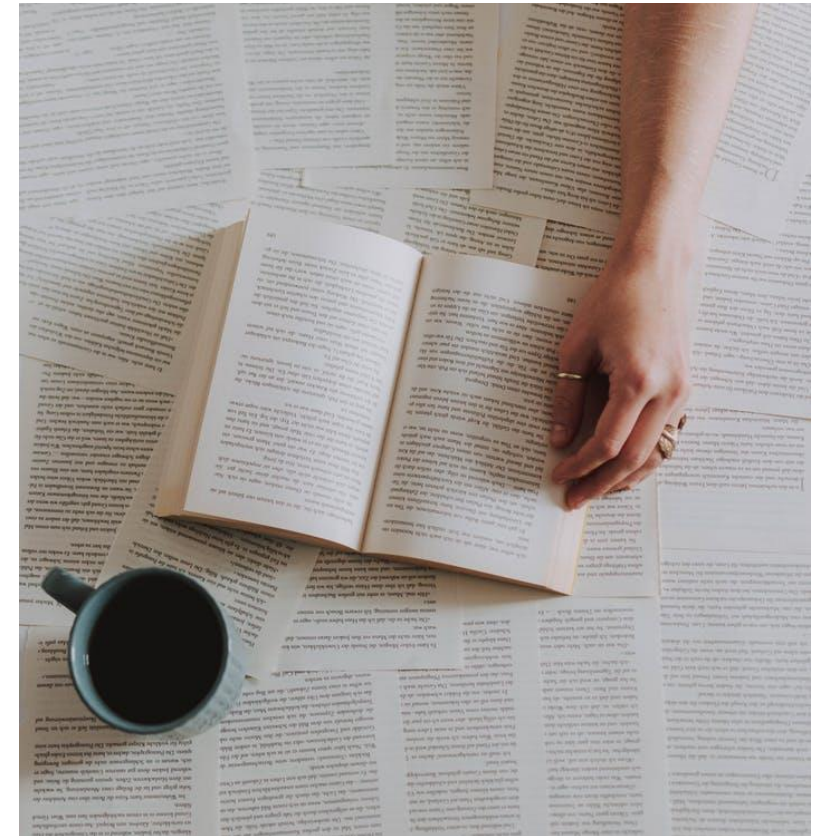- Split up into groups of 3-6 students. Proposal midway, check-ins.

**Final Project**:
- 40% total, groups of 3-6; proposal midway. **More info soon!**

# **Class Setup**: Reading

No textbooks
- I will post useful notes, primers, papers
  - See schedule page

- Expect **new papers** (submitted during the timeframe of the class)

- For presentations: we will have a list of papers to pick from, but new/unlisted papers are options as well

# **Class Setup**: Background

More on this at the end of class, but

- **Basic ML** (at the level of 760 or so)
  - Short review next lecture
- **Technical components:**
  - Linear Algebra
  - Calculus
  - Probability

Note: this class is partially **conceptual** and partially **technical**

# **Class Setup**: Goals

Two goals:

- Become acquainted with **how to use** large pretrained/language/foundation models
- Understanding the technical underpinnings of these models and **_why_** they work

**Note**: if you are only interested in a very broad overview of ML, then CS 540 or 760 might be a better choice.

# **Class Setup**: Goals II

Mini-goals:
- **Understanding** research

- **Big picture/**ML ecosystem

- **Intuition** around modern ML paradigms

# Break & Questions

# What We'll Cover

- The past: **supervised learning**

  - **Dataset:**
  $$\left(x^{(1)}, y^{(1)}\right), \left(x^{(2)}, y^{(2)}\right), \ldots, \left(x^{(n)}, y^{(n)}\right)$$

  
  safe

  
  poisonous

  
  safe

  - **Model:**

  
  Simonyan and Zisserman

  **Train:**

  

# New Paradigms: Pretraining

# New Architectures: Transformers

- Lots of architectures (RNNs, LSTMs, CNNs) → Transformers



Vaswani et al. '17

# New Models

A bunch you've seen already,



And more all the time:

# New Interfaces: Prompting



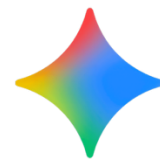**FR** Explain how transformers work

Transformers are a type of deep learning model that has gained significant popularity and success in various natural language processing (NLP) and machine learning tasks. They are the foundation for many state-of-the-art models, including BERT, GPT, and more. Here's a high-level explanation of how transformers work:

1. **Self-Attention Mechanism:**

   At the core of the transformer architecture is the self-attention mechanism. This mechanism allows the model to weigh the importance of different parts of the input data, giving it the ability to consider context from all positions simultaneously. The self-attention mechanism works as follows:

   - **Query, Key, and Value Vectors:** For each input token or word, three vectors are computed: Query, Key, and Value. These vectors are learned during the training process and are used to capture relationships between words.

   - **Attention Scores:** The model computes a compatibility score (dot

**Prompt**

University professor clearly explaining machine learning to a class

# Reasoning

- Chain-of-thought and friends:



**Standard Prompting**

Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

Model Output

A: The answer is 27. ✖
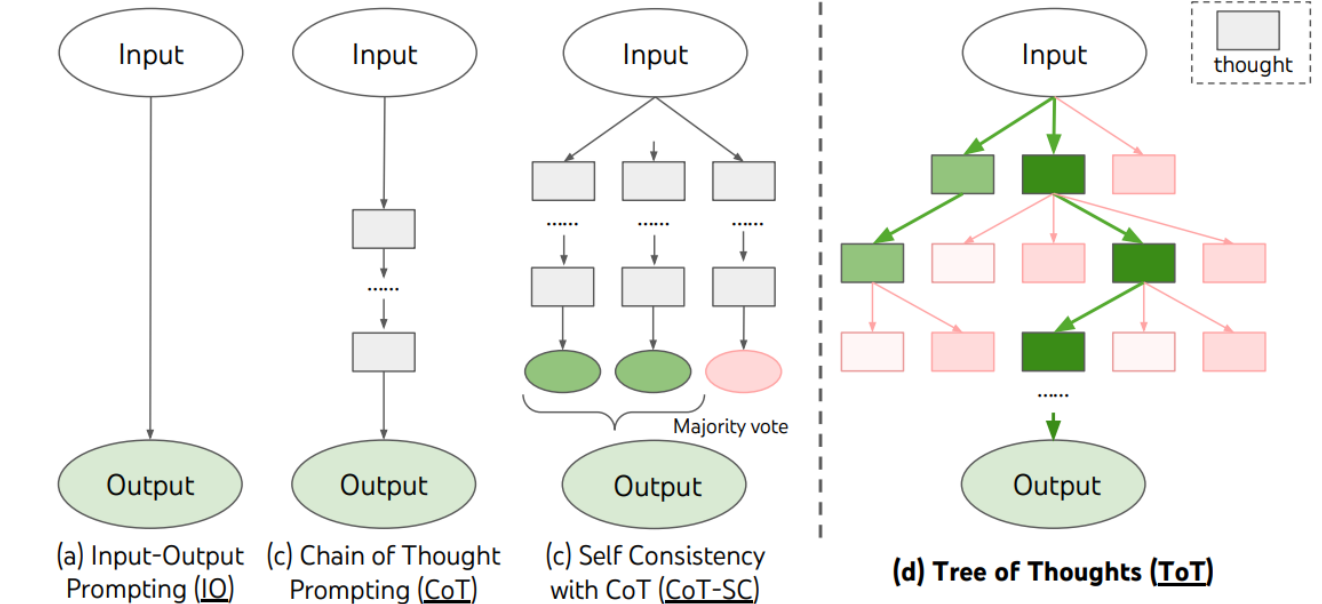
Wei et al

**Chain-of-Thought Prompting**

Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls do...

A: Roger starte... each is 6 tenni...

Q: The cafeter... make lunch an... do they have?

Model Outpu...

A: The cafeteri... 20 to make lun... bought 6 more... answer is 9. ✔

(a) Input-Output Prompting (IO)

(c) Chain of Thought Prompting (CoT)

(c) Self Consistency with CoT (CoT-SC)

Majority vote

(d) Tree of Thoughts (ToT)

thought

Yao et al

# Adapting & Improving Models

- Prompt Engineering
- Fine-tuning
- Steering
- RL-based tuning



Cuenca and Paul

# Model Alignment

- RLHF, DPO, and more!



Low quality data

Text
e.g. Internet data

Optimized for
text completion

Language
modeling

Pretrained LLM

High quality data

Demonstration
data

Finetuned for
dialogue

Supervised
finetuning

SFT model

Human feedback
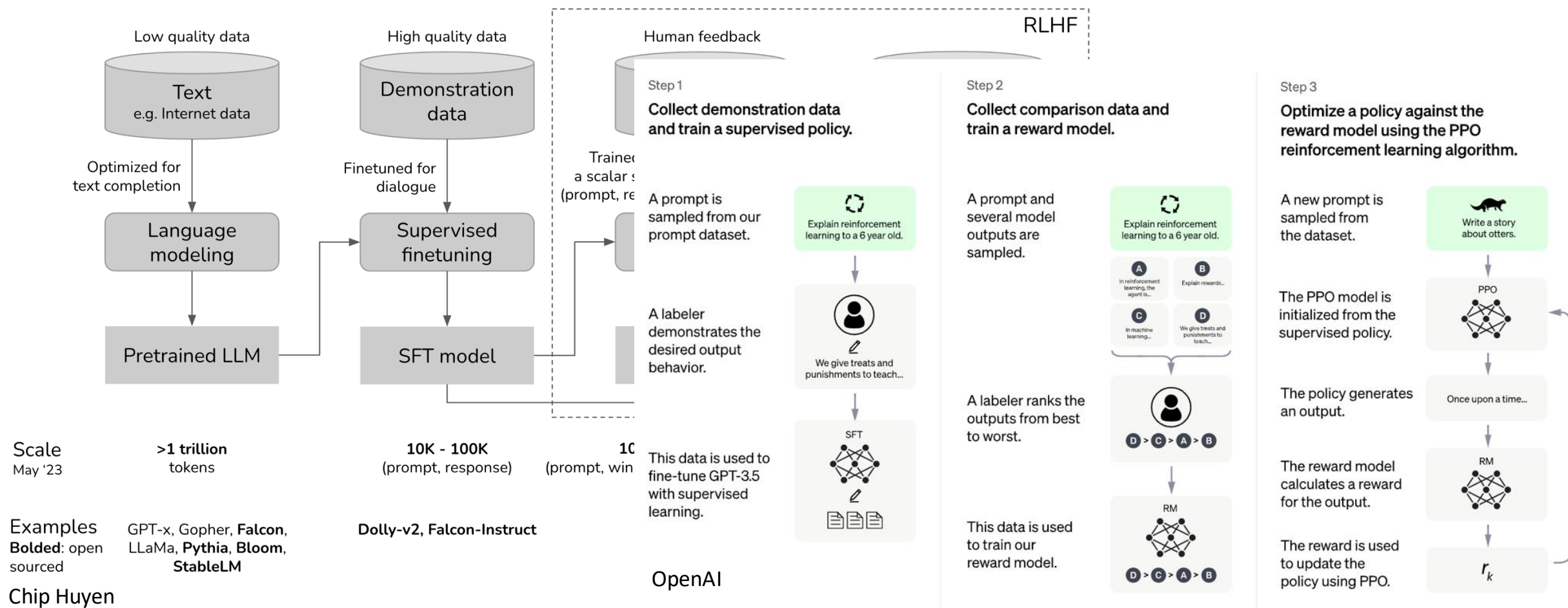
Trained
a scalar s
(prompt, re

RLHF

**Step 1**
**Collect demonstration data**
**and train a supervised policy.**

A prompt is
sampled from our
prompt dataset.

Explain reinforcement
learning to a 6 year old.

A labeler
demonstrates the
desired output
behavior.

We give treats and
punishments to teach...

SFT

This data is used to
fine-tune GPT-3.5
with supervised
learning.

**Step 2**
**Collect comparison data and**
**train a reward model.**

A prompt and
several model
outputs are
sampled.

Explain reinforcement
learning to a 6 year old.

A  B
In reinforcement
learning, the
agent is...   Explain rewards...
C  D
In machine
learning...   We give treats and
punishments to
teach...

A labeler ranks the
outputs from best
to worst.

D > C > A > B

RM

This data is used
to train our
reward model.

D > C > A > B

**Step 3**
**Optimize a policy against the**
**reward model using the PPO**
**reinforcement learning algorithm.**

A new prompt is
sampled from
the dataset.

Write a story
about otters.

The PPO model is
initialized from the
supervised policy.

PPO

The policy generates
an output.

Once upon a time...

The reward model
calculates a reward
for the output.

RM

The reward is used
to update the
policy using PPO.

$r_k$

Scale
May '23

Examples
**Bolded**: open
sourced

Chip Huyen

>1 trillion
tokens

GPT-x, Gopher, **Falcon**,
LLaMa, **Pythia**, **Bloom**,
**StableLM**

10K - 100K
(prompt, response)

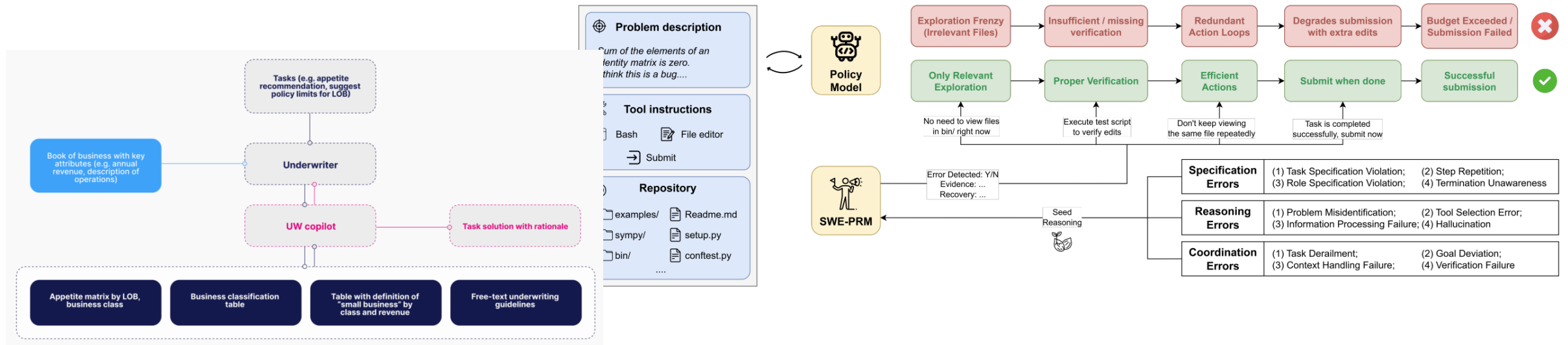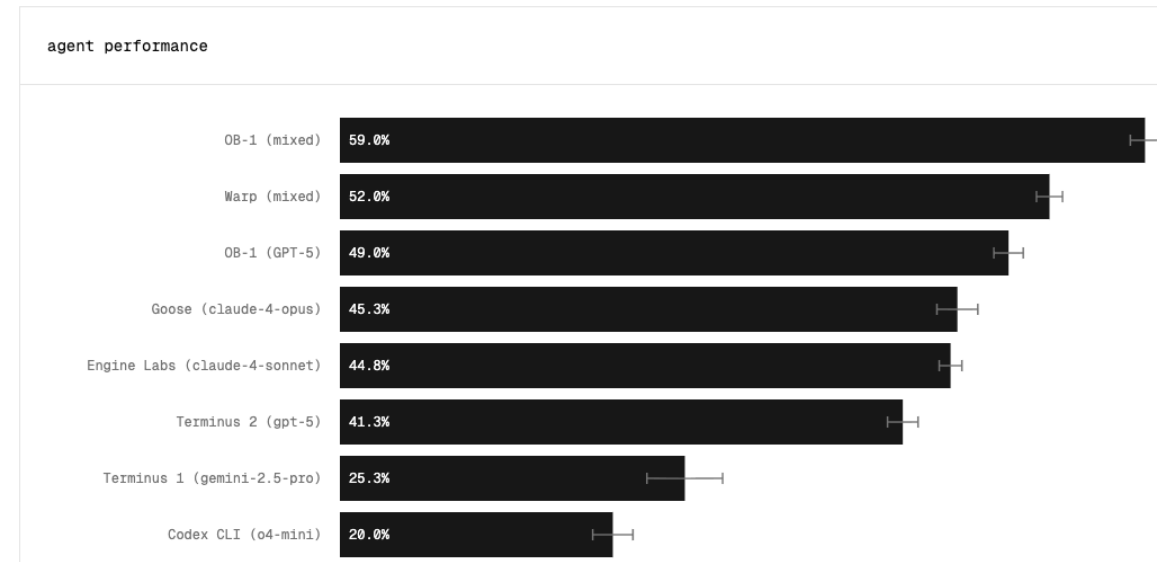**Dolly-v2, Falcon-Instruct**

10
(prompt, win

OpenAI

# Evaluating Models

- Benchmarks
- LLM-as-a-judge



Huang et al '24

# Reasoning, RL, Agents

- Reasoning capabilities
- GRPO
- Autonomous agents
- Environments

# Training & Data

# Multimodal Models



LAVR
- Training Paradigms (§3.4.2)
  - Training-free
    - Few-shot: Chameleon [75], HuggingGPT [80], MM-REACT [14], VISPROG [76]
    - Zero-shot: CAT [84], CaFo [90], Visual ChatGPT [83], IdealGPT [91], ChatCaptioner [92], PointCLIP V2 [93], SMs [94]
  - Finetuning: GPT4Tools [34]
- Functions (§3.4.3)
  - Controller: HuggingGPT [80], GPT4Tools [34], Chameleon [75], VISPROG [76]
  - Decision Maker: Visual ChatGPT [83], ChatCaptioner [92], MM-REACT [14], IdealGPT [91]
  - Semantics Refiner: CAT [84], PointCLIP V2 [93], CaFo [90], SMs [94]
- Evaluation (§3.4.4)
  - Benchmark Metrics: Chameleon [75], IdealGPT [91], PointCLIP V2 [93], CaFo [90], VISPROG [76], SMs [94]
  - Manual Evaluation: ChatCaptioner [92], VISPROG [76], GPT4Tools [34]

Yin et al

# Scaling

Scaling laws:

# Security, Privacy, Bias

Some of the issues we'll encounter...

# THE DARK SIDE OF LARGE LANGUAGE MODELS

Part 2: "Who's a good chatbot?"

By: Eoin Wickens, Marta Janus

# Break & Questions

# Brief History of Foundation Models

**Three Historical Trends**
- Brief introduction, more to come.

1. **Multitask** models (old!)
2. **Pretraining** and fine-tuning (~2015-)
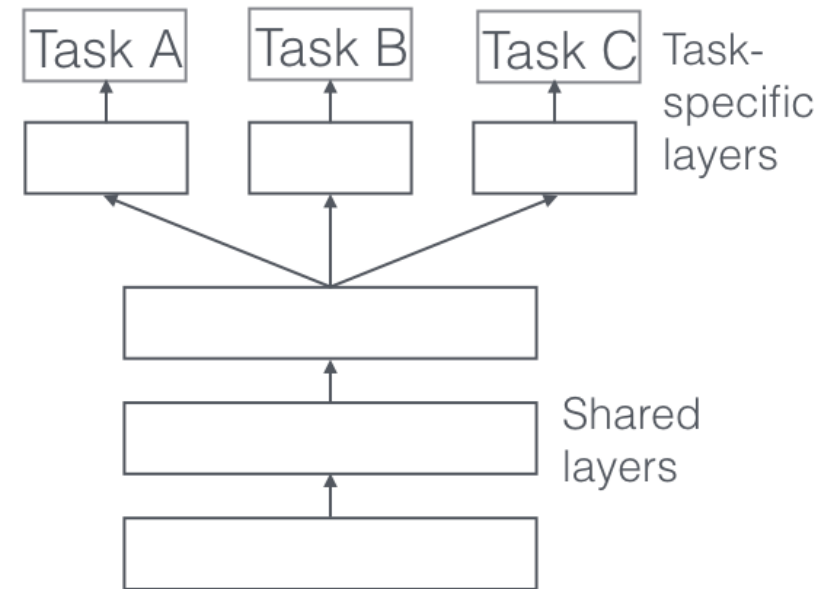3. **Word embeddings** and language models (~2013-)

# 1. Multitask Models

**Idea:** Given tasks $T_1$, ..., $T_k$, rather than training k separate models, train a common base and task-specific "heads"

- Related to ***transfer learning***

Differences (vs. modern FMs)

- Usually fixed tasks
- Train on data from all tasks (limited)

J. Ray

# 2. Pretraining and Fine-tuning

**Motivation:** Training from scratch is expensive. Why?

- Deep learning revolution (2010-). Confluence of
  - Larger datasets (ImageNet etc)
  - Larger hardware resources (GPU support)
  - Produces larger models

- Much of 2010-2015 CV research builds larger and larger CNNs, so training costs ↑
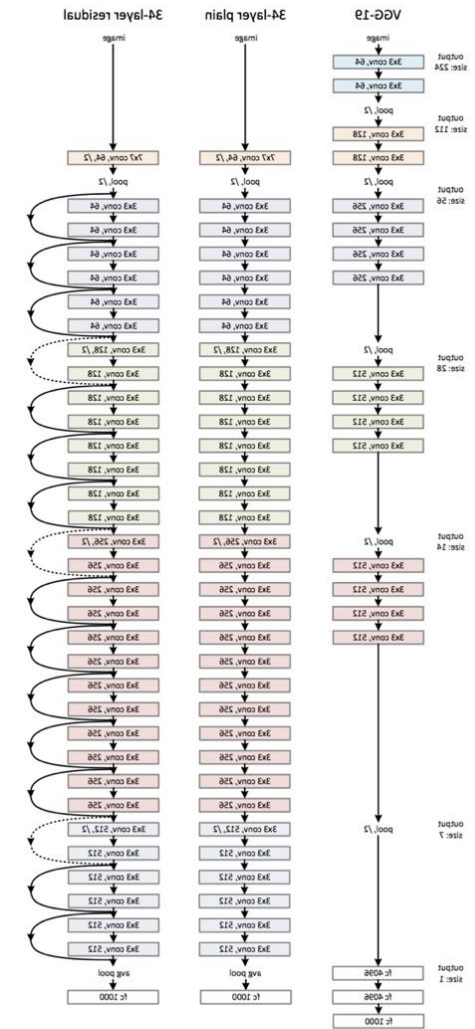
Figure 3. Example network architectures for ImageNet. **Left:** the VGG-19 model [41] (19.6 billion FLOPs) as a reference. **Middle:** a plain network with 34 parameter layers (3.6 billion FLOPs). **Right:** a residual network with 34 parameter layers (3.6 billion FLOPs). The dotted shortcuts increase dimensions. Table 1 shows more details and other variants.

He et al '16.

# 2. Pretraining and Fine-tuning

**Motivation:** Training from scratch is expensive.

**Idea**: *pretrain* a single model on a dataset
- Then *fine-tune* to adapt to downstream task
- Ex: pretrained ResNets on ImageNet (2015-)

**Issues**:
- Other data modalities/domains? Could build ImageNet analogue, but expensive
- Leads to **self-supervised training** (2016-)
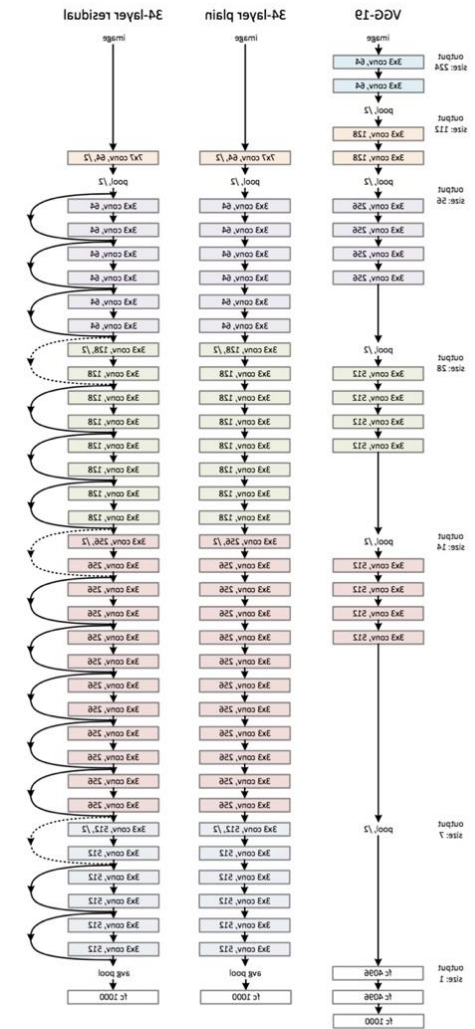  - No labels needed! Ex: SimCLR, DINO, lots more



Figure 3. Example network architectures for ImageNet. **Left**: the VGG-19 model [41] (19.6 billion FLOPs) as a reference. **Middle**: a plain network with 34 parameter layers (3.6 billion FLOPs). **Right**: a residual network with 34 parameter layers (3.6 billion FLOPs). The dotted shortcuts increase dimensions. **Table 1** shows more details and other variants.
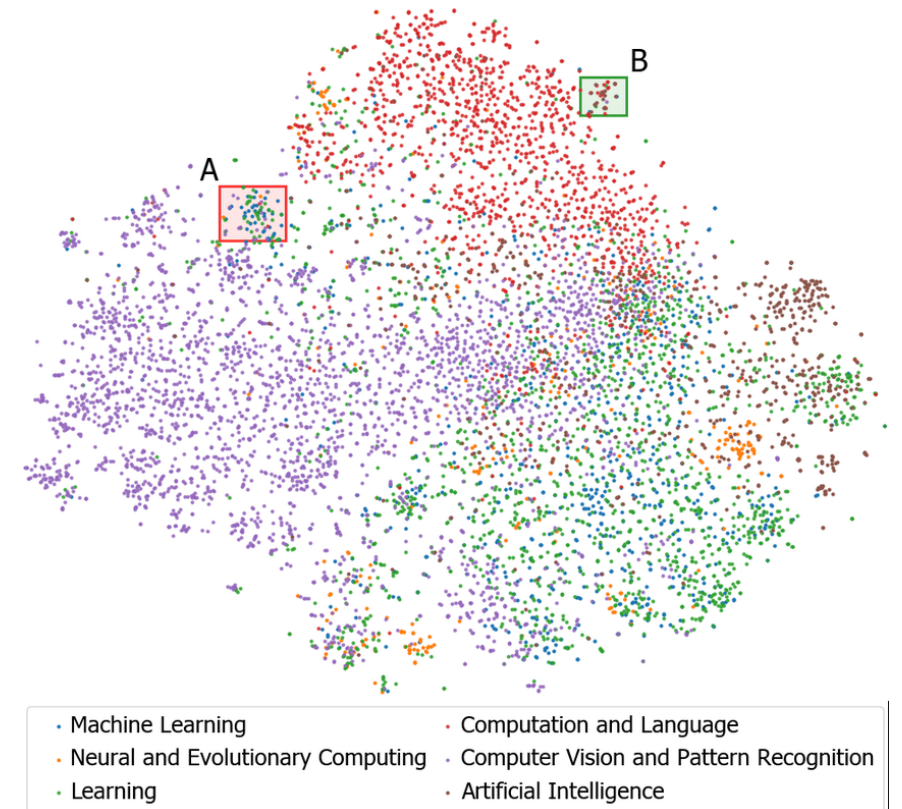
He et al '16.

# 3. Word Embeddings and Language Models

**Motivation:** Deep learning advances –
can they be applied to NLP?

Three areas of application:

1. General: ***word embeddings***
2. Specific: ***translation tasks***
3. Specific: ***language modeling tasks***

Embeddings for arXiv papers (6 ML categories)

- Machine Learning
- Neural and Evolutionary Computing
- Learning
- Computation and Language
- Computer Vision and Pattern Recognition
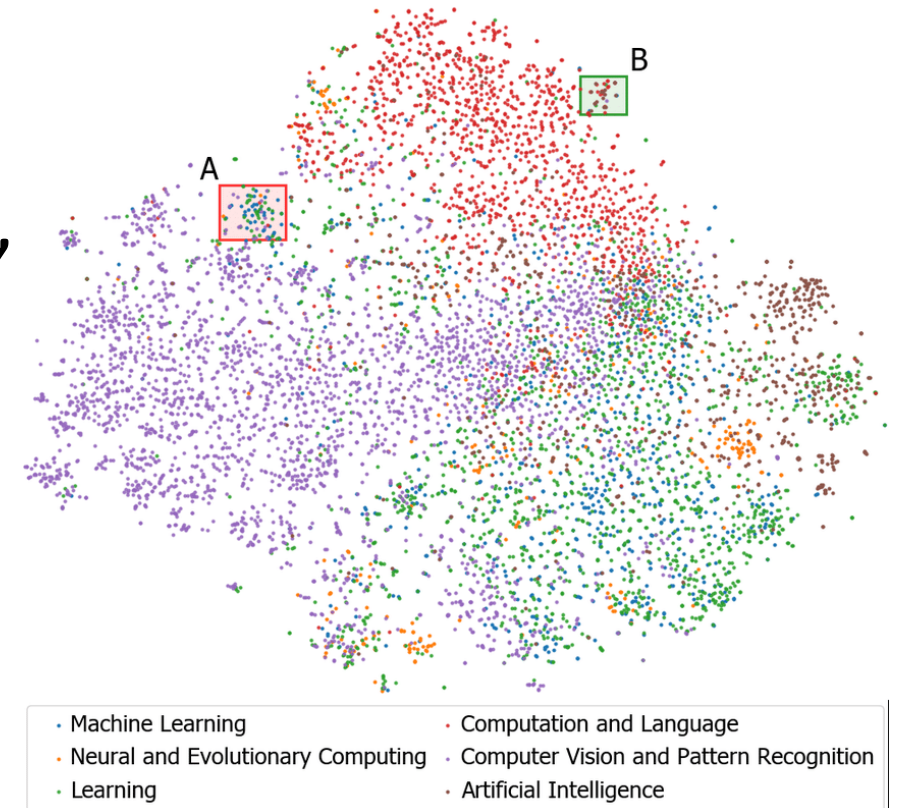- Artificial Intelligence

Lo et al '19.

# 3. Word Embeddings and Language Models

**Motivation:** Can we learn, in advance, structured representations of words?

- Then plug into language-specific networks (LSTMs, etc)?

- Word embeddings (2013-2016): Glove, Word2Vec, etc.
  - A form of *representation learning*

- **Issues:** static. No context used for words like "bank"

Embeddings for arXiv papers (6 ML categories)

- Machine Learning
- Neural and Evolutionary Computing
- Learning
- Computation and Language
- Computer Vision and Pattern Recognition
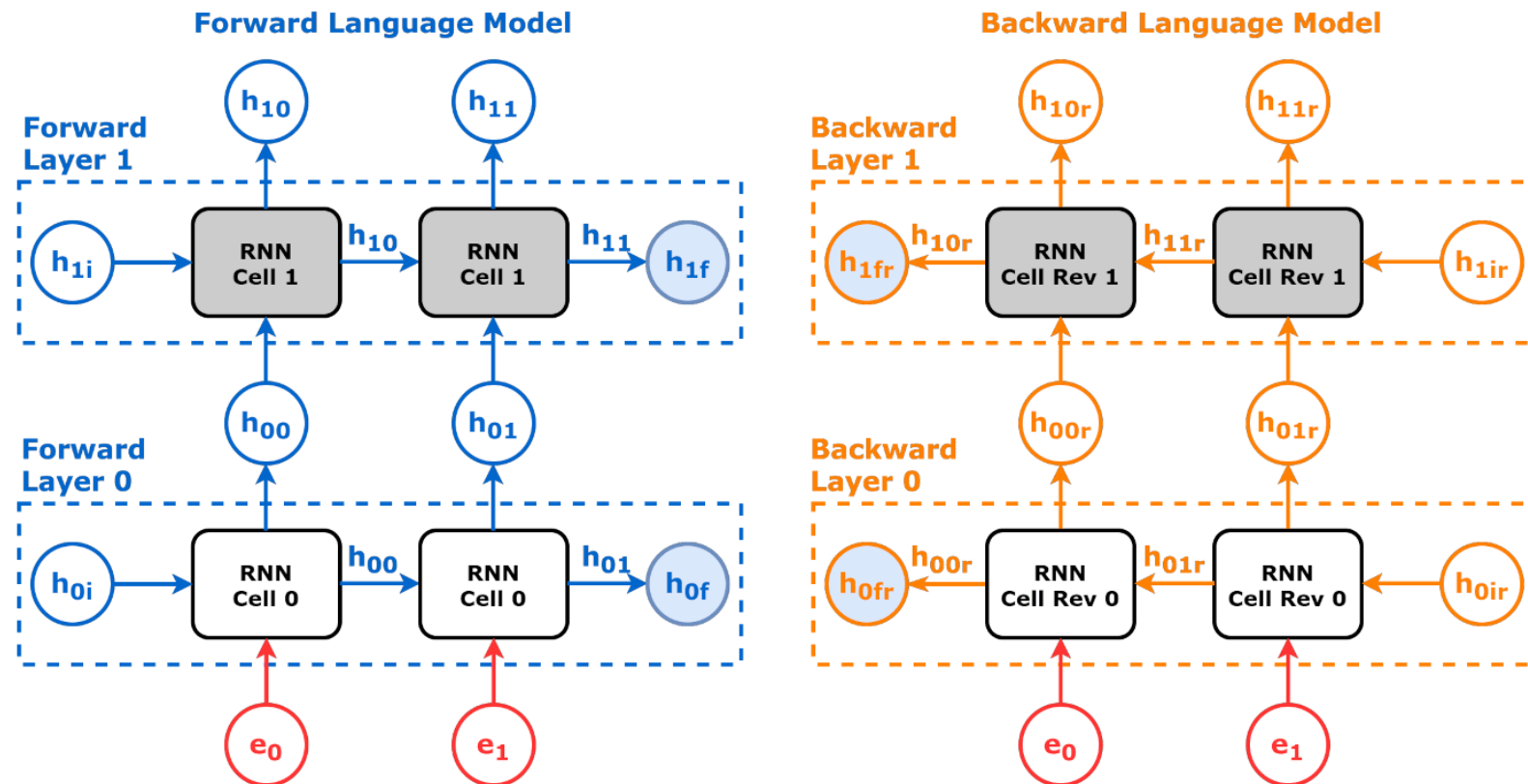- Artificial Intelligence

Lo et al '19.

# 3. Word Embeddings and Language Models

**Solution:** Contextual word embeddings

- **Idea:** Plug into a model to obtain the embedding, and include the context
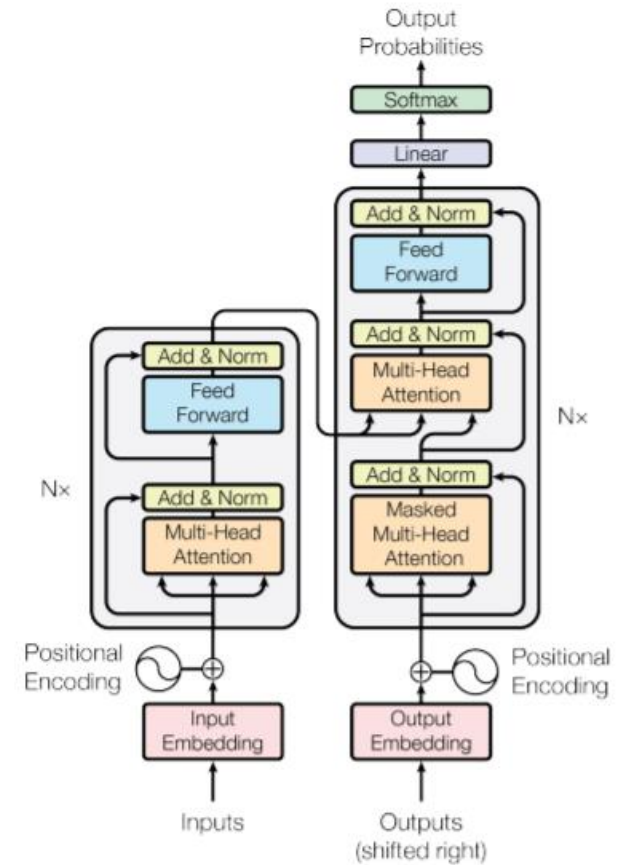
- **ELMO embeddings:**



Godoy

# 3. Word Embeddings and Language Models

**So far:** embeddings, which are general (whether static or contextual)

- What about deep learning advances for specific tasks?

**Translation:** critical task

- New architecture: ***Transformers*** (2017)
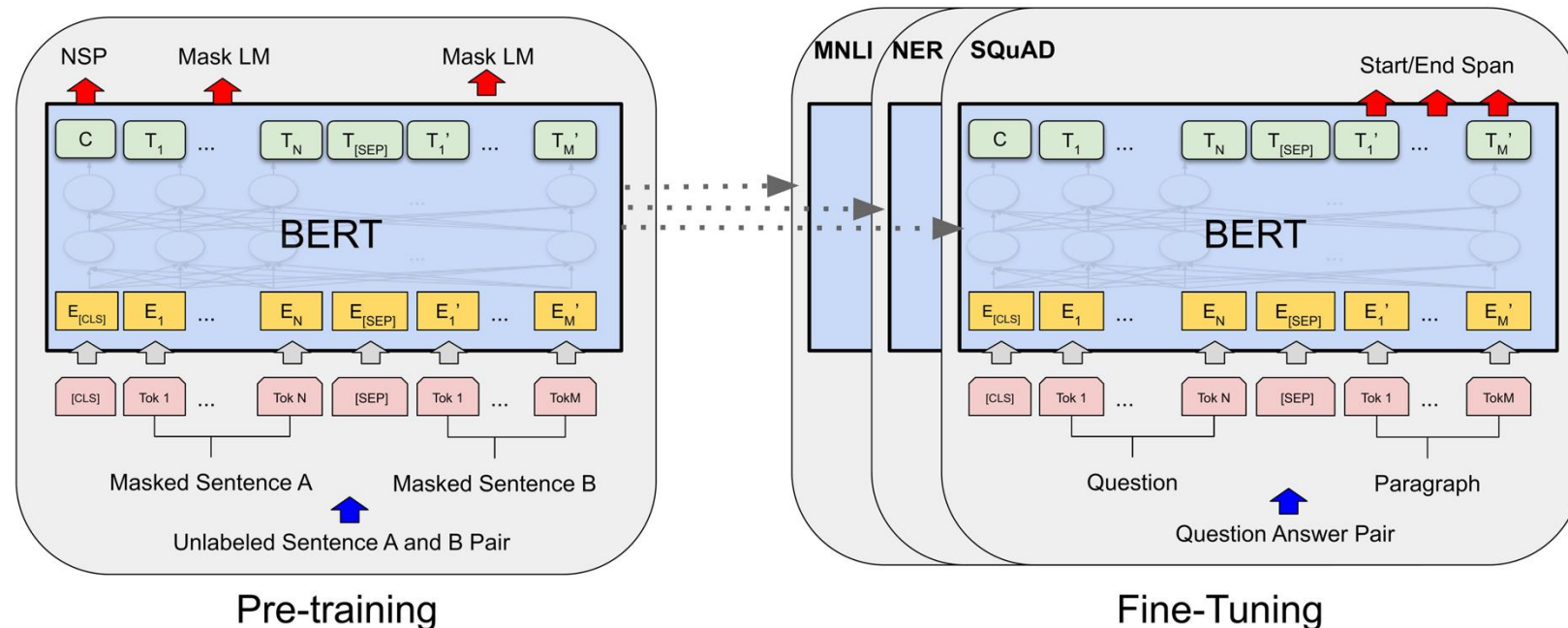- Uses ideas around attention (2014-)



Vaswani et al. '17

# 3. Word Embeddings and Language Models

**So far:**

- Contextual embeddings (ELMO)
- Translation via Transformers architecture

Combine to ***BERT***, perhaps the first modern foundation model



Devlin et al. '18

# 3. Word Embeddings and Language Models

**What about language models?**

- Similar idea: replace older architecture language models with new Transformer architecture

- Ex: **GPT** (**G**enerative **P**retrained **T**ransformer)

- In all cases, pretrain on massive text corpora
  - All the way back to static embeddings, use all of Wikipedia!

# Summary

**Modern foundation models**

- Build on old ideas about multitask learning,
- Are large-scale and pretrained on massive data, then specialized
  - Dating back to vision models from mid 2010s
- First heavily scaled for NLP applications, building on ideas on
  - Powerful contextual word embeddings
  - New architectures suitable for text (and beyond)

# Next two weeks

1.  **Review of ML**
    - Very short!
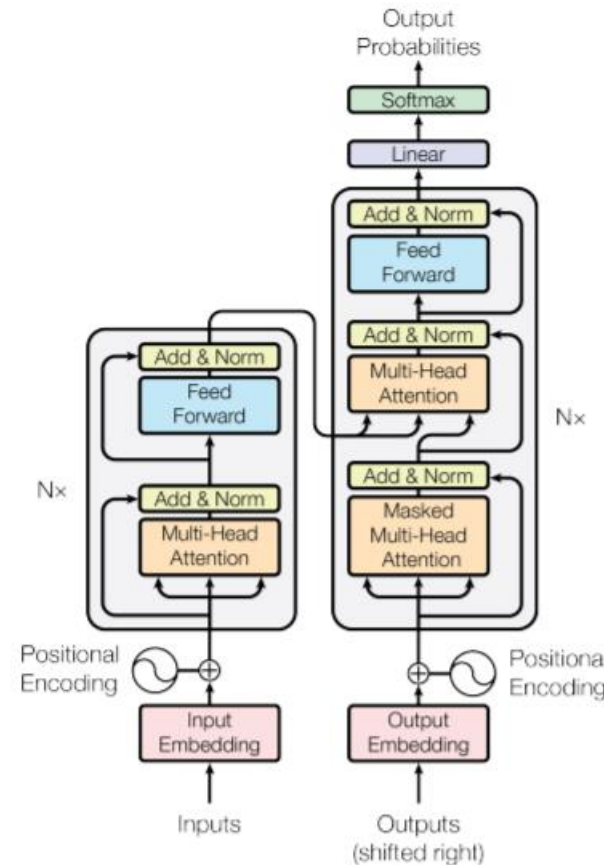
2.  **Architectures:**
    1.  **Transformers**
        - Intro to attention
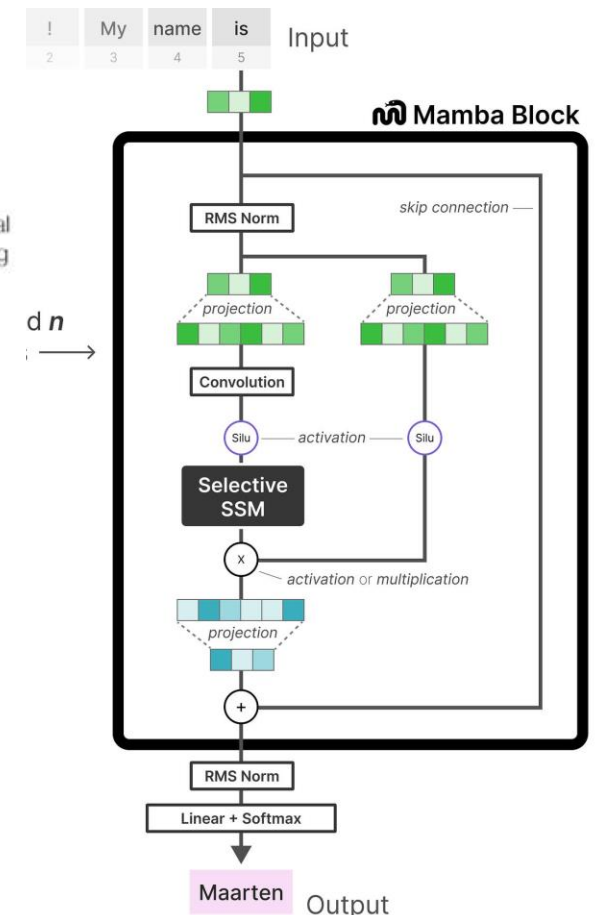    2.  **Sub-quadratic architectures**

3.  **Language Models**
    - Encoder-decoder, Encoder-only, Decoder-only, etc

Vaswani et al. '17

M. Grootendorst