# AI & LLM Overview

1. Clarify some terms

2. Transformers

3. Large Language Models

4. Pre-training vs Fine-tuning

# LLMs in AI Safety

1. Clarify more terms

2. What is AI safety?

3. Examples of technical research

4. Resources

Google Developer Groups

```
ext(
 'Section Title',
 style: TextStyle(
   color: Colors.green[200],
  ),
),
```
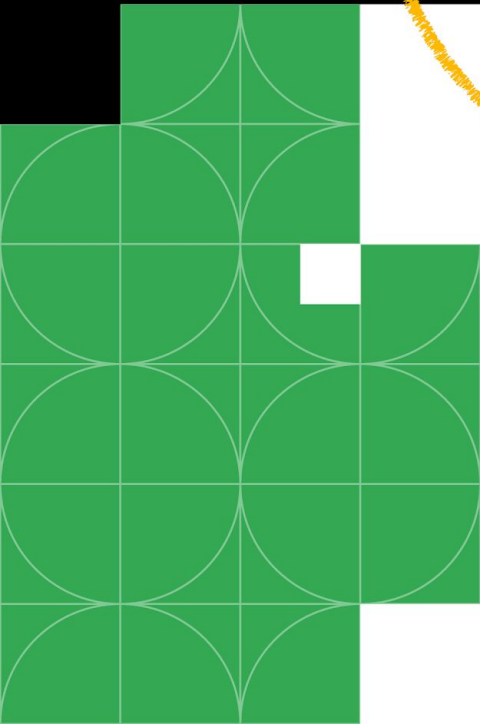
```
s.star,
r: Colors.green[500],

Text('23'),
```

**devfest**

Google Developer Groups

Sacramento

# AI & LLM Overview

# Terminology

**AI**: Artificial Intelligence

**ML**: Machine Learning

**LLM**: Large Language Model

**NLP**: Natural Language Processing

**GPT**: Generative Pretrained Transformer

**Transformers**: Neural network leading to LLMs

**RLHF**: Reinforcement Learning from Human Feedback

# Natural Language Processing [NLP]:
## Computers can speak & understand human languages

Pre-1990s:
**Rule-based Expert Systems**

# 1990s-2000s:
# Statistics & Probabilities

You shall know a word by the company it keeps

J.R. Firth, Linguist
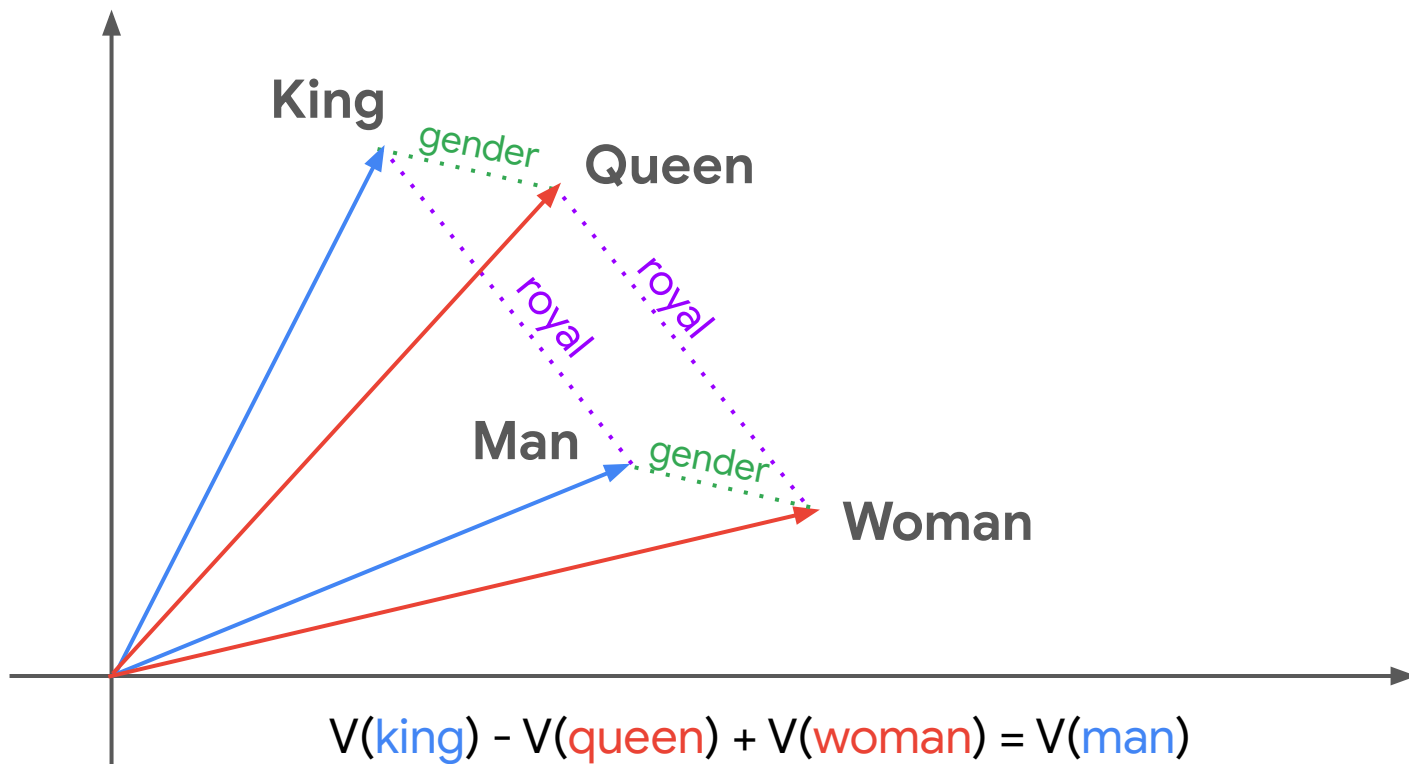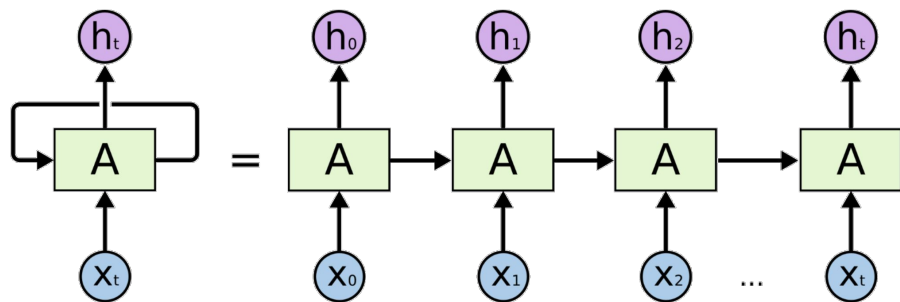
# 2010s:
## Rise of Deep Learning and Neural Networks

2013: Word2Vec Embeddings

V(king) - V(queen) + V(woman) = V(man)

# 2013: Word2Vec Embeddings

| Analogies | Word Pair 1 | | Word Pair 2 | |
|---|---|---|---|---|
| Man-Woman | king | queen | man | woman |
| Capital city | Athens | Greece | Oslo | Norway |
| City-in-state | Chicago | Illinois | Sacramento | California |
| Opposite | possibly | impossibly | ethical | unethical |
| Comparative | great | greater | tough | tougher |
| Nationality adjective | Switzerland | Swiss | Canada | Canadian |
| Past tense | walking | walked | swimming | swam |
| Plural nouns | mouse | mice | dollar | dollars |

# 2010s:
# Neural Networks
RNN, GRU, LSTM



Google Developer Groups

# Early Neural Networks

- Slow & forgetful

# 2017: Transformers

- Attention mechanism
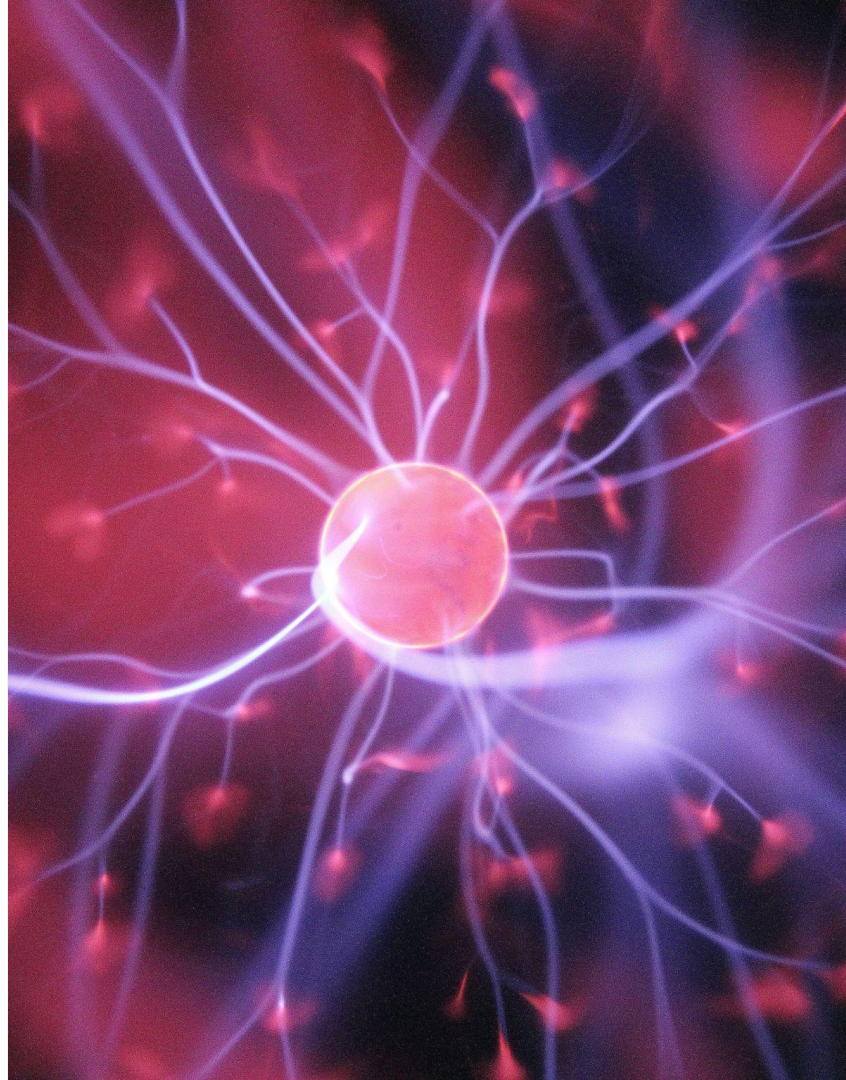- Parallel processing
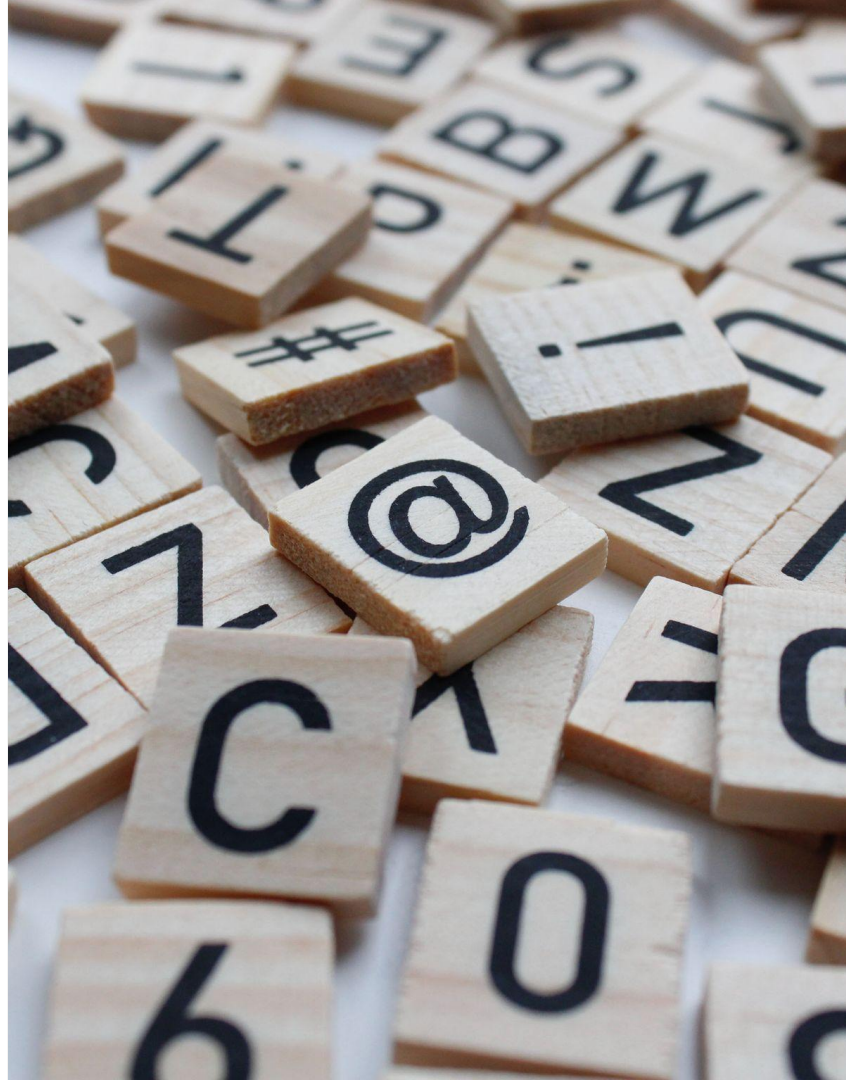
# Transformer Architecture
## Encoder + Decoder



Google Developer Groups

# Rise of LLMs

>1 billion neurons



Google Developer Groups

Trained for
**next word prediction**

# Pre-trained base

vs

# Fine-tuned models

Google Developer Groups

**RLHF:**
Reinforcement Learning
from Human Feedback

Google Developer Groups

# Fine-tuned to
**follow instructions**

# Fine-tuned for
# **conversations**

# What's Next?

- Multimodal

- Open-source

- Agents

devfest

Google Developer Groups

Sacramento

# LLMs in
# AI Safety

# More Terminology

**AGI**: Artificial GENERAL Intelligence, Strong AI

**ANI**: Artificial NARROW Intelligence, Weak AI

**HLI**: Human-Level Intelligence

**ASI**: Artificial SUPERINTELLIGENCE surpass humans on all tasks

**Orthogonality Thesis**: intelligence & goals are independent

**Takeoff**: fast vs slow

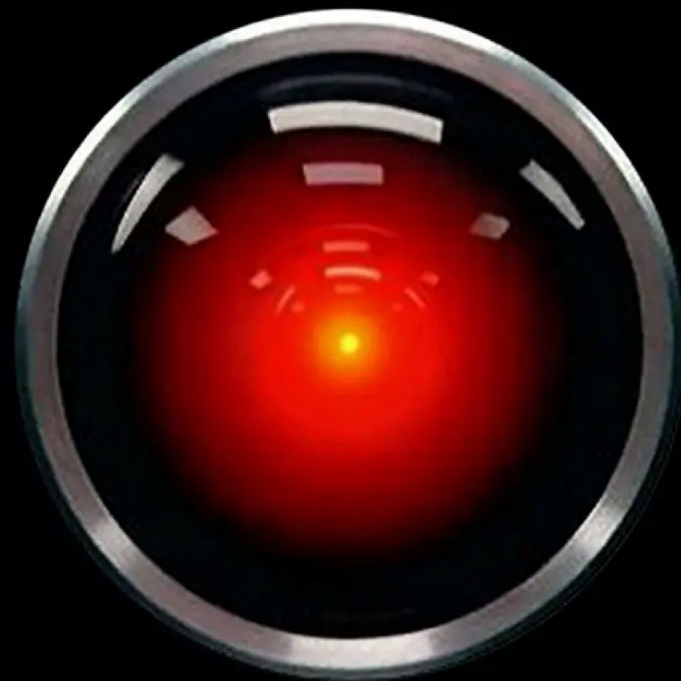**Timelines**: short vs long

# ANI

Artificial Narrow Intelligence
Weak AI

# AGI

Artificial General Intelligence
Strong AI

# Intelligence:

- Human-level
- Superintelligence
   ... Singularity

# Takeoff
Fast vs Slow
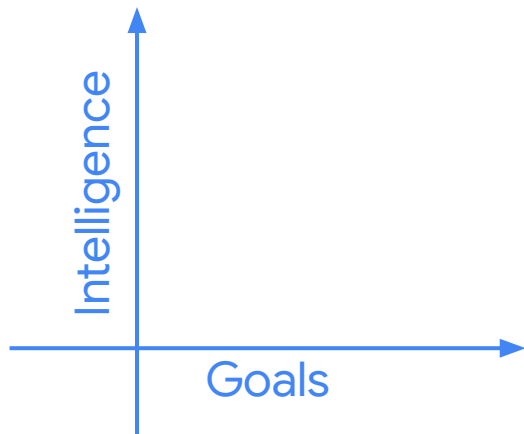
# Timelines
Short vs Long

# Orthogonality Thesis

Intelligence & Goals are independent



Google Developer Groups

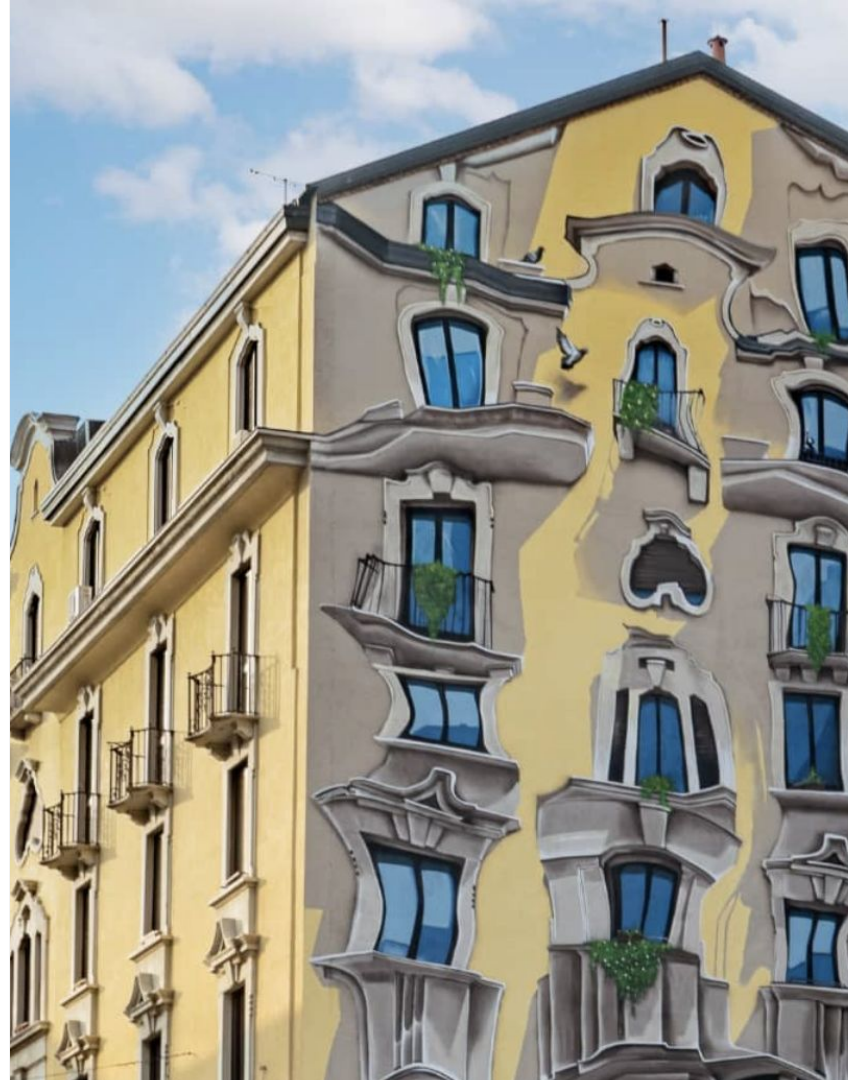# Paperclip Maximizer

Thought experiment
or stamp alternative

# Outer Misalignment
AI creator goals don't align with general human values

# Inner Misalignment
AI achieves goals in ways unintended by its creators

Google Developer Groups

# Gorilla Problem

Can humans maintain autonomy in a world with superintelligence?

Google Developer Groups

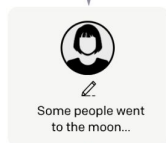# Some Concrete Examples of Technical Research...

# Alignment: RLHF Human Values

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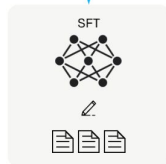**Collect demonstration data, and train a supervised policy.**

A prompt is sampled from our prompt dataset.

Explain the moon landing to a 6 year old

A labeler demonstrates the desired output behavior.
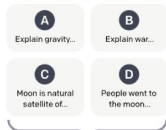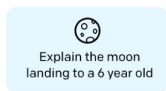
Some people went to the moon...

This data is used to fine-tune GPT-3 with supervised learning.

SFT

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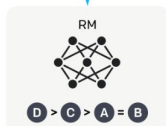**Collect comparison data, and train a reward model.**

A prompt and several model outputs are sampled.

Explain the moon landing to a 6 year old

A — Explain gravity...
B — Explain war...
C — Moon is natural satellite of...
D — People went to the moon...

A labeler ranks the outputs from best to worst.

D > C > A = B

This data is used to train our reward model.

RM

D > C > A = B

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

A new prompt is sampled from the dataset.

Write a story about frogs

The policy generates an output.

PPO

Once upon a time...

The reward model calculates a reward for the output.

RM

The reward is used to update the policy using PPO.

$r_k$

Google Developer Groups

# **LLM Output Evaluation:** Inverse Scaling Laws

# LLM Internals
Interpretability

Google Developer Groups

# Resources

1. **AI Explained videos on AI development + safety**
   youtube.com/@aiexplained-official

2. **80,000 Hours career advice + job board**
   80000hours.org/problem-profiles/artificial-intelligence

3. **AISafety.info FAQs**
   AISafety.info

4. **AI Safety Fundamentals online curricula**
   AISafetyFundamentals.com

5. **Alignment Forum share research + discussions**
   AlignmentForum.org

Google Developer Groups

# Questions?

Embrace
Safely

# devfest

## ChengCheng Tan

https://www.linkedin.com/in/cheng2-tan/

ccstan99@gmail.com

 Google Developer Groups

Sacramento