# Human in the Loop

**Rishabh Anand**

**@rishabh16_**

# The GPT Series

## GPT → Generative Pretrained Transformer
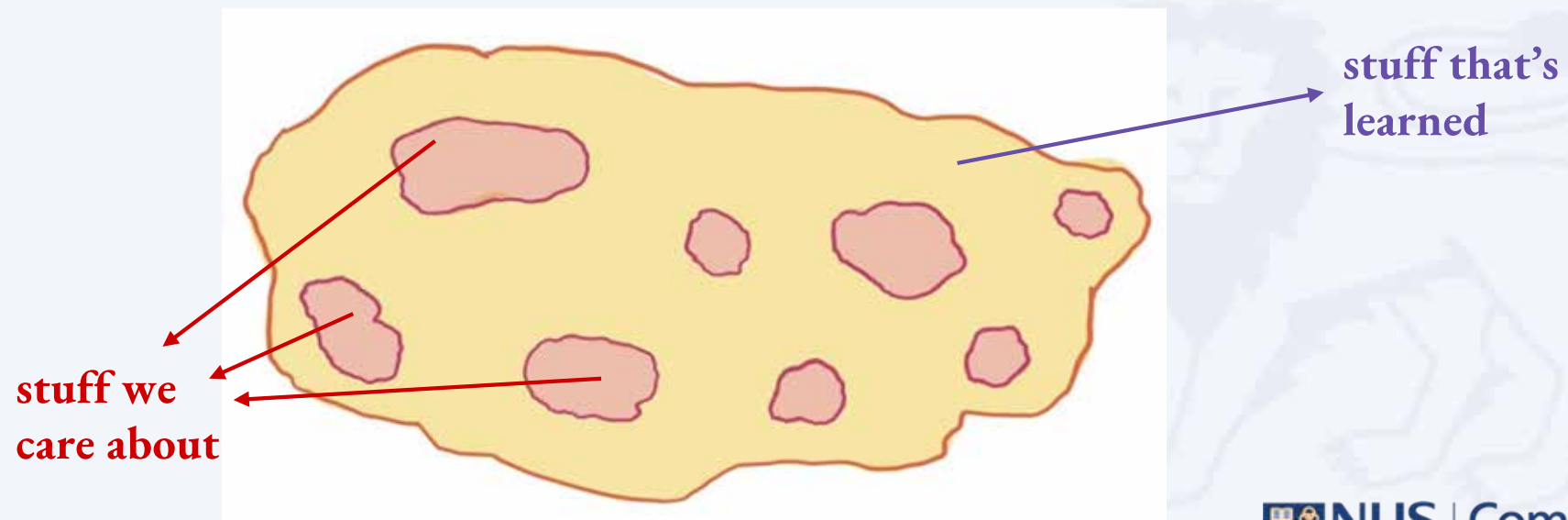
**Rapid growth ...**

| 2018 | 2019 | 2020 | 2022 | 2023 |
|------|------|------|------|------|
| GPT | GPT-2 | GPT-3 | **ChatGPT**(GPT-3.5) | **GPT-4** |
| *Improving Language Understanding by Generative Pre-Training* | *Language Models are Unsupervised Multi-task Learners* | *Language Models are Few-Shot Learners* | ***Training language models to follow instructions with human feedback*** * | (technical report) |

**\* GPT-3.5 is built on top of *InstructGPT* with a different data collection setup**

**NUS | Computing**
National University

# Reinforcement Learning from Human Feedback

# (Large) Language Models

- **Language Models (like GPT-X),**
  - are chaotic
  - model a "giant mass of people" ~ Minqi Jiang, MetaAI

- **For different prompts, you can get wildly different outputs**

- **We must "snip out" the ugly, less-preferred parts**

# RL from Human Feedback

- **Provides a friendlier interface to interact with LMs**

- **Biases the underlying model to generate human-aligned content**

- **Improves reliability, honesty, and safety of LLMs**

**"How do we get LLMs to sound more *human*?"**

# RL from Human Feedback

# RL from Human Feedback

1. **Pretrain a LLM on a body of text [GPT-X, for instance]**

# RL from Human Feedback

1. **Pretrain a LLM on a body of text [GPT-X, for instance]**

2. **Train a Reward Model (RM) → "how would a human feel?"**

# RL from Human Feedback

1. Pretrain a LLM on a body of text [GPT-X, for instance]

2. Train a Reward Model (RM) → "how would a human feel?"

3. Finetune using RL [LLM agent predicts words and is scored]

# LLMs + RLHF



[source]

# ChatGPT
# for
# Students

# Ask Away!

- **Treat ChatGPT as you would a friend**

- **Want something? Just ask for it!**

- **The art of "Prompt Engineering" with ChatGPT**

**Use ChatGPT as a personal tutor!**

# The Possibilities

- **Digestible explanations**

- **Summarising Long-form content**

- **Peer Review + feedback**

# Generate Digestible Explanations

# Summarising Content

**Given some long-form content that contains a lot to go through ...**

# Summarising Content

# Peer Review + Feedback

# Peer Review + Feedback

Please edit the following machine learning paper to make it clearer: "In this section, we propose a model for shifting seasonal matrix factorization. A tensor we consider consists of a timestamped series of (m × n) matrices X(1), . . . , X(r), which can be sparse, until the current time point r. We incrementally observe a new matrix X(r + 1) and r evolves (r = r + 1).

Our goal is to forecast X(t) where r < t by uncovering important factors in the flow of data, whose characteristics can change over time. As we discussed the effectiveness of handling seasonal patterns, we incorporate seasonal factors into a switching model that can adaptively recognize recent patterns in forecasting."

"In this section, we propose a new model for analyzing data that changes over time and is affected by seasonal patterns. Our model considers a series of matrices, X(1), X(2), ..., X(r), that are labeled with a time stamp and may have missing information. As new matrices, X(r+1), come in, we update our analysis. Our goal is to predict future matrices, X(t), where t is a time after r, by finding key patterns in the data that change over time and taking into account patterns that repeat in certain seasons. We use a flexible model that adapts to recent patterns to improve our predictions."

# Peer Review + Feedback

# ChatGPT for Students

- **LLM technology will only get better from here on**

- **Students ~~should~~ *can* learn how to operate these tools**

- **While LLMs can improve productivity, it's not the be-all-end-all**

**AI tools lower the <u>activation energy</u> to get started!!!**

# But ... shortcomings?

**Stay for our panels!**