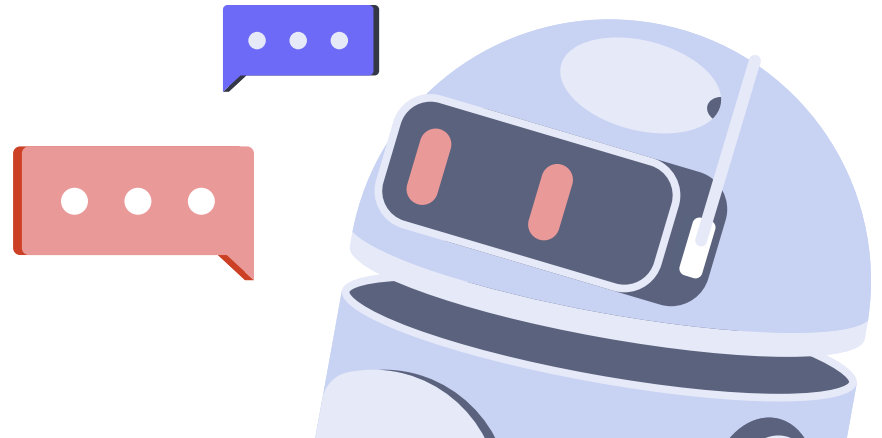# Boosting Business or a Risky Gamble?

Enterprise Perspectives on Large Language Models

Szilvia Hodvogner
2023 ML Forum

# Large Language Models (LLMs)

2018
GPT

2019
GPT2

2020
GPT3

2022
GPT3.5,
Flan-T5

2023
GPT4, LLaMa,
Alpaca, Dolly, …

# Large Language Models (LLMs)

| 2018 | 2019 | 2020 | 2022 | 2023 |
|------|------|------|------|------|
| GPT | GPT2 | GPT3 | GPT3.5, Flan-T5 | GPT4, LLaMa, Alpaca, Dolly, … |

- Transformers with attention
- Billions of parameters
- Petabytes of data from the internet
- Advancements:
  - N-shot learning
  - Reinforcement Learning from Human Feedback (RLHF)
  - Instruction-tuning
  - …

# Large Language Models (LLMs)

| 2018 | 2019 | 2020 | 2022 | 2023 |
|------|------|------|------|------|
| GPT | GPT2 | GPT3 | GPT3.5, Flan-T5 | GPT4, LLaMa, Alpaca, Dolly, … |

- **Transformers** with attention
- **Billions** of parameters
- **Petabytes of data** from the internet
- Advancements:
  - N-shot learning
  - Reinforcement Learning from Human Feedback (RLHF)
  - Instruction-tuning
  - …

**Goal: Predict the next word → Limitless potential!**

# Limitless potential... BUT

**Ethical and social challenges**

Model-level constraints

Technology-level limitations

**Responsible AI Regularization**

# Limitless potential… BUT

**Ethical and social challenges**

↓

**Responsible AI Regularization**

- **Disruptive**
- **Effects? - jobs, teaching, content creating**
- **Legal policies? - training data, biases, alignment**

# Limitless potential... BUT

| Ethical and social challenges | Model-level constraints | Technology-level limitations |
|---|---|---|

Responsible AI Regularization

Parameters, configurations, legal questions

# Closed vs open-source models

| | Closed development | | Open-source model | |
|---|---|---|---|---|
| | **GPT3.5** | **GPT4** | **LLaMa** | **Dolly** |
| # params | 175B | 1000B | 65B | 12B |
| Context token size | 4K | 32K | 8K | 2K |
| Inference RAM needs | ~750GB | Unknown | ~130GB | ~32GB |
| Pricing ($/1K tokens) | 0.002 | 0.03/0.06 | Free | Free |
| Commonsense reasoning (HellaSwag dataset) | 86% | 95% | 84% | 71% |
| Q&A, pronoun resolution (WinoGrande dataset) | 82% | 88% | 77% | 62% |

**Params**

**Resources**

**Results**

# Closed vs open-source models

| | Closed development | | Open-source model | |
| --- | --- | --- | --- | --- |
| | **GPT3.5** | **GPT4** | **LLaMa** | **Dolly** |
| # params | 175B | 1000B | 65B | 12B |
| Context token size | 4K | 32K | 8K | 2K |
| Inference RAM needs | ~750GB | Unknown | ~130GB | ~32GB |
| Pricing ($/1K tokens) | 0.002 | 0.03/0.06 | Free | Free |
| Commonsense reasoning (HellaSwag dataset) | 86% | 95% | 84% | 71% |
| Q&A, pronoun resolution (WinoGrande dataset) | 82% | 88% | 77% | 62% |

**Params**

**Resources**

**Results**

# Closed vs open-source models

| | Closed development | | Open-source model | |
| --- | --- | --- | --- | --- |
| | **GPT3.5** | **GPT4** | **LLaMa** | **Dolly** |
| # params | 175B | 1000B | 65B | 12B |
| Context token size | 4K | 32K | 8K | 2K |
| Inference RAM needs | ~750GB | Unknown | ~130GB | ~32GB |
| Pricing ($/1K tokens) | 0.002 | 0.03/0.06 | Free | Free |
| Commonsense reasoning (HellaSwag dataset) | 86% | 95% | 84% | 71% |
| Q&A, pronoun resolution (WinoGrande dataset) | 82% | 88% | 77% | 62% |

**Params**

**Resources**

**Results**

Percentage of correct answers

# Who owns what?

| Closed development | | Open-source model | |
|---|---|---|---|
| **GPT3.5** | **GPT4** | **LLaMa** | **Dolly** |
| **Commercial use** | **Commercial use** | **Non-commercial use** | **Commercial use** |

# Who owns what?

| Closed development | | Open-source model | |
| --- | --- | --- | --- |
| **GPT3.5** | **GPT4** | **LLaMa** | **Dolly** |
| **Commercial use** | **Commercial use** | **Non-commercial use** | **Commercial use** |

## Not simple!

- Chinchilla: closed development, non-commercial use

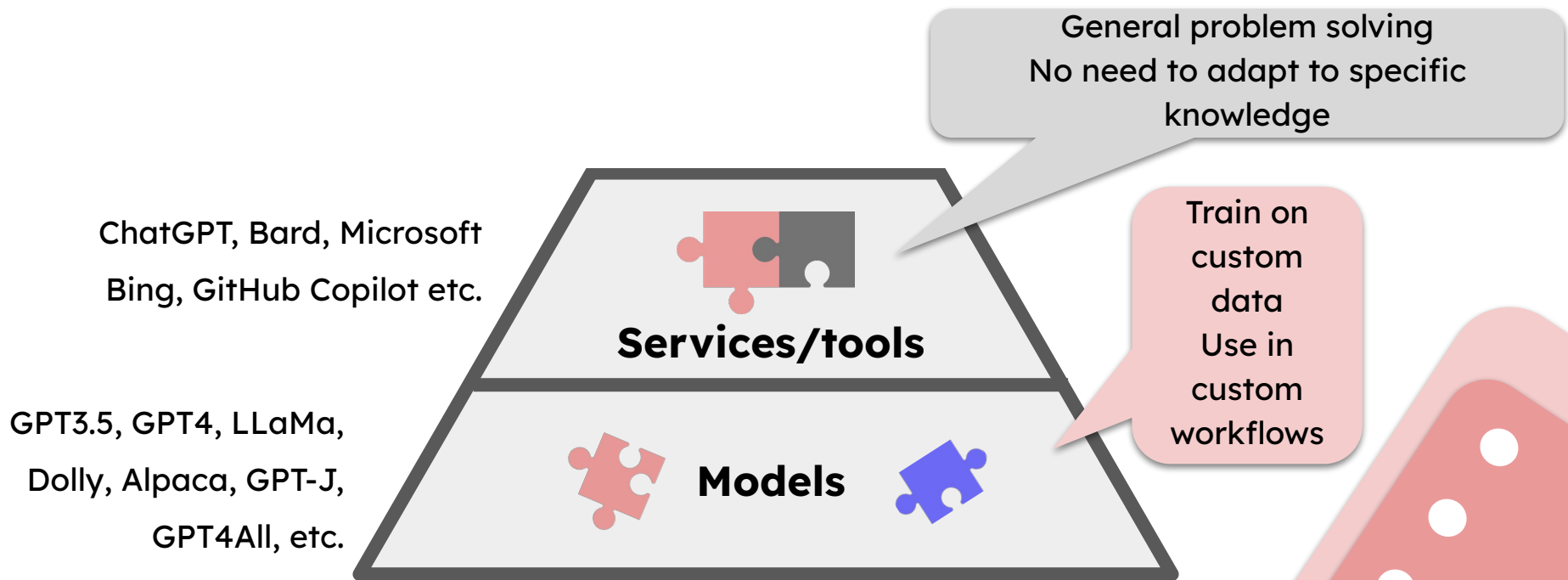- Vicuna: open-source, commercial use licence, but still non-commercial use because LLaMa

# Building blocks
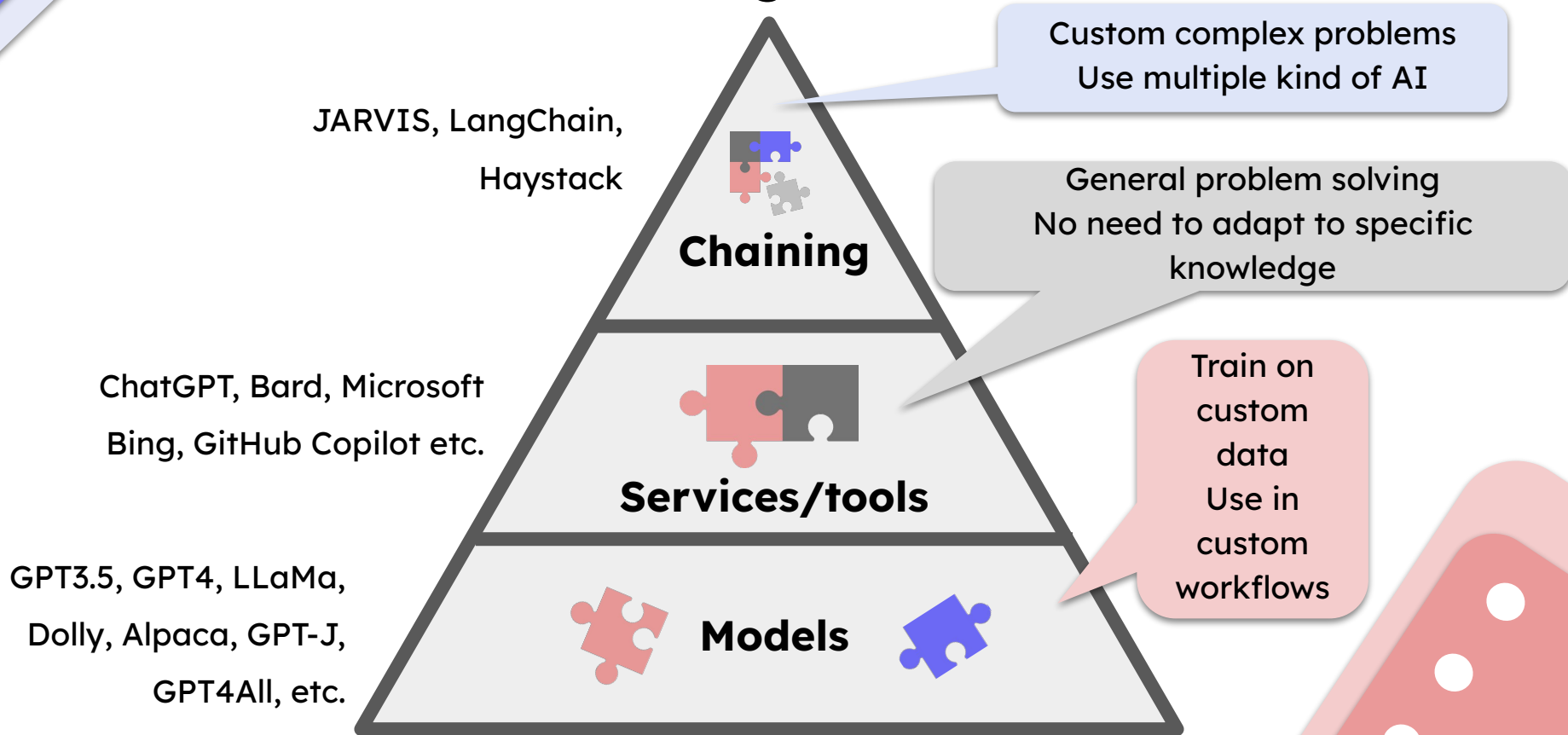
GPT3.5, GPT4, LLaMa,
Dolly, Alpaca, GPT-J,
GPT4All, etc.

**Models**

Train on custom data
Use in custom workflows

# Building blocks

ChatGPT, Bard, Microsoft
Bing, GitHub Copilot etc.

GPT3.5, GPT4, LLaMa,
Dolly, Alpaca, GPT-J,
GPT4All, etc.

**Services/tools**

**Models**

General problem solving
No need to adapt to specific
knowledge

Train on
custom
data
Use in
custom
workflows

# Building blocks

JARVIS, LangChain, Haystack

**Chaining**

Custom complex problems
Use multiple kind of AI

General problem solving
No need to adapt to specific knowledge

ChatGPT, Bard, Microsoft Bing, GitHub Copilot etc.

**Services/tools**

Train on custom data
Use in custom workflows

GPT3.5, GPT4, LLaMa, Dolly, Alpaca, GPT-J, GPT4All, etc.

**Models**

# Limitless potential... BUT

Ethical and social
challenges

Model-level
constraints

Technology-level
limitations

Responsible AI
Regularization

Parameters,
configurations,
legal questions

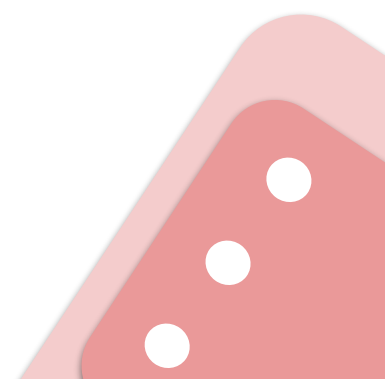Factfullness,
security, privacy

# Factual incorrectness and hallucinations

**Hallucinations:**

A confident response, but it is nonsensical or unfaithful to the provided source content

# Factual incorrectness and hallucinations

**Hallucinations:**

A confident response, but it is nonsensical or unfaithful to the provided source content

**Cause:**

Predicting the next word, no reasoning capabilities

# Factual incorrectness and hallucinations

A confident response, but it is nonsensical or unfaithful to the provided source content

**Cause:**
Predicting the next word, no reasoning capabilities

**Mitigation strategies:**
- More data and training
- Fine-tune
- Teaching how to say "I don't know"
- Middle tier: Query databases, search internet, other models, knowledge graphs

# Security & Privacy

Closed-developments - What happens to the upload data?
- **Privacy** - Data goes through
- **Security** - Rely on the API

**OpenAI**
*"Retain it 30 day, than dispose"*

**Azure OpenAI**
*"Opt out, they cannot use your data"*

**ChatGPT**
*"Remove chat history"*

# Security & Privacy

Closed-developments - What happens to the upload data?
- **Privacy** - Data goes through
- **Security** - Rely on the API

Do not share:
- Sensitive data about company, person, or customer
- Proprietary codes
- Emails, meetings drafts

Always check and validate AI response!

**OpenAI**
*"Retain it 30 day, than dispose"*

**Azure OpenAI**
*"Opt out, they cannot use your data"*

**ChatGPT**
*"Remove chat history"*

# Security & Privacy

Closed-developments - What happens to the upload data?
- Privacy - Data goes through
- Security - Rely on the API

Open-source solutions
- Privacy - You are safe
- Security - Your concern

**OpenAI**
*"Retain it 30 day, than dispose"*

**Azure OpenAI**
*"Opt out, they cannot use your data"*

**ChatGPT**
*"Remove chat history"*

# Security & Privacy

Closed-developments - What happens to the upload data?
- Privacy - Data goes through
- Security - Rely on the API

Open-source solutions
- Privacy - You are safe
- Security - Your concern

Chaining or combining LLMs with tools
- Caution with autonomous actions!
  - Harmful queries, drop database, send email to wrong contacts, insert obfuscated code, call unintended APIs

**OpenAI**
*"Retain it 30 day, than dispose"*

**Azure OpenAI**
*"Opt out, they cannot use your data"*

**ChatGPT**
*"Remove chat history"*

# Integrating LLMs into a company

| | **API calls** - closed devs<br><br>(e.g. chatGPT, GPT3.5) | **Local hosting** - open devs<br><br>(e.g. Dolly, LLaMa, Alpaca) |
|---|---|---|
| **For** ✔️ | • Easy setup<br>• Cutting edge models are available<br>• Low latency | • Higher selection of models<br>• Lower inference cost<br>• Independent |
| **Against** ✖️ | • Privacy concerns<br>• Higher long-term costs<br>• Depends on a 3rd party API | • Complex and costly setup<br>• Weaker models |

# Future

| Ethical and social challenges | Model-level constraints | Technology-level limitations |
|---|---|---|
| ↓ | ↓ | ↓ |
| Learn to beneficially integrate into life | Smaller models Better open-source alternatives | Chaining Plugins Sparks of AGI? |

# Thank you for your attention!

Enterprise Perspectives on Large Language Models