# Large Language Models in finance

Laurent Sorber

Financieel Forum Vlaams-Brabant
Leuven, 20 November 2023

InvestSuite

Radix

# Outline

- Who are InvestSuite & Radix
- What are Large Language Models
- The magic of LLMs
- Dispelling the magic of LLMs
- Applications of LLMs
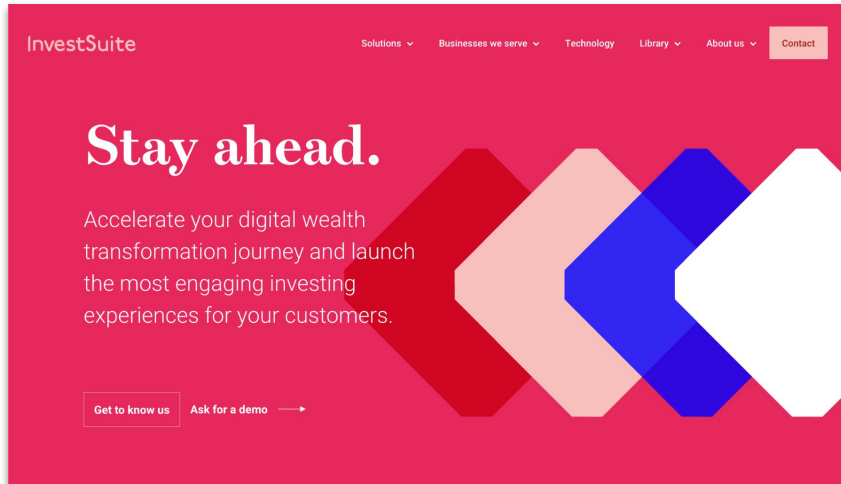- What does the future hold

# Outline

- **Who are InvestSuite & Radix**
- What are Large Language Models
- The magic of LLMs
- Dispelling the magic of LLMs
- Applications of LLMs
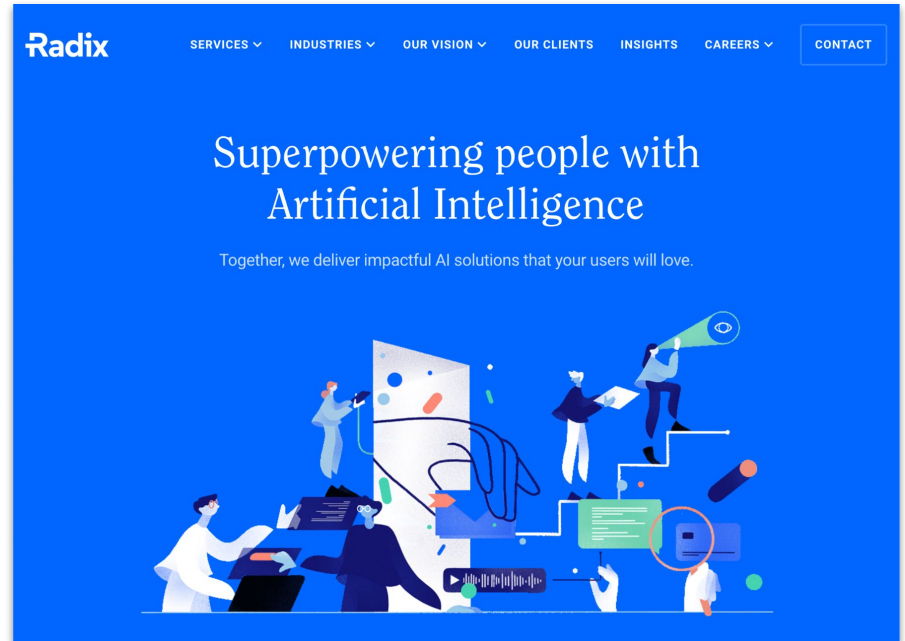- What does the future hold

# InvestSuite

**A suite of investtech products** with which we want to make investing more accessible to a wider audience



# Radix

**Bespoke AI solutions** with which we want to superpower people for clients across industries and AI domains

# InvestSuite

**A suite of investtech products** with which we want to make investing more accessible to a wider audience

Founded 2018

Offices in Leuven, CH, UK

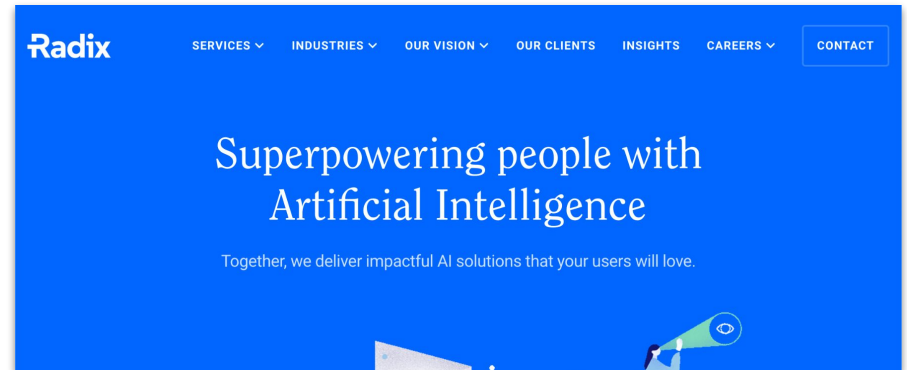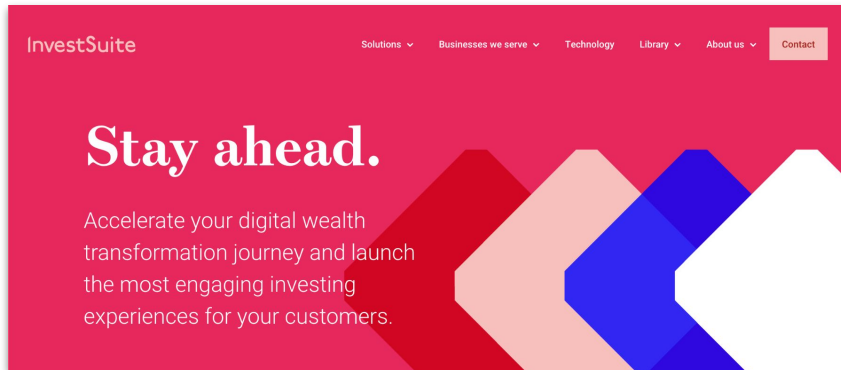50+ colleagues

Sales presence in BE, DE, NL, USA

# Radix

**Bespoke AI solutions** with which we want to superpower people for clients across industries and AI domains
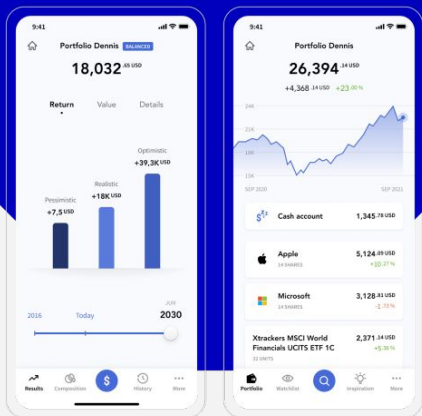
Founded 2018

Offices in Brussels & Ghent

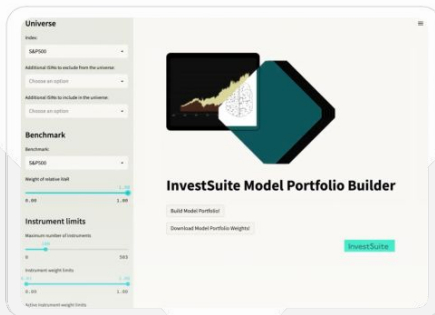50+ colleagues
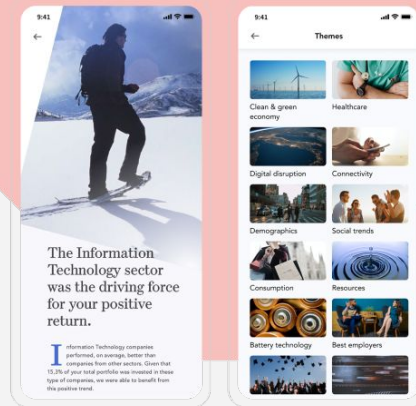
Active in BE, NL, DE, IT, Nordics, USA

---

InvestSuite

Solutions ⌄    Businesses we serve ⌄    Technology    Library ⌄    About us ⌄    Contact

## Stay ahead.

Accelerate your digital wealth transformation journey and launch the most engaging investing experiences for your customers.

---

Radix

SERVICES ⌄    INDUSTRIES ⌄    OUR VISION ⌄    OUR CLIENTS    INSIGHTS    CAREERS ⌄    CONTACT

## Superpowering people with Artificial Intelligence

Together, we deliver impactful AI solutions that your users will love.

**InvestSuite**

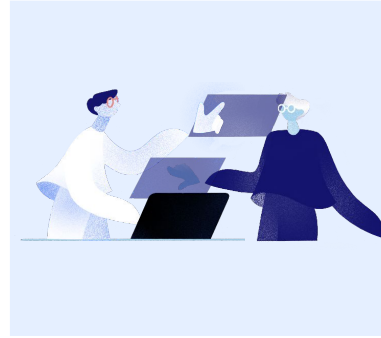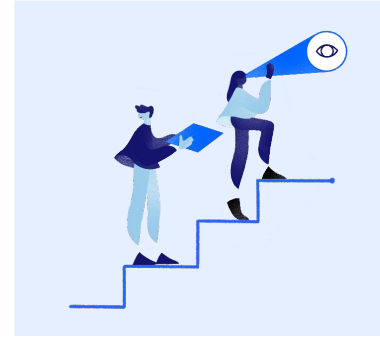| Investment Platforms | Portfolio Construction | Investment Insights & Reporting |
|---|---|---|
| Robo Advisor    Self Investor | Portfolio Optimizer | StoryTeller    Insight API |

**Radix**

## AI solutions

- AI modelling
- Software engineering
- Cloud infrastructure
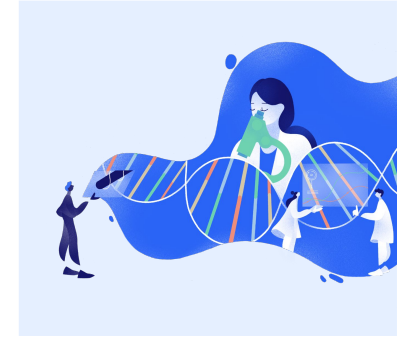- MLOps
- Data annotation

## Knowledge transfer

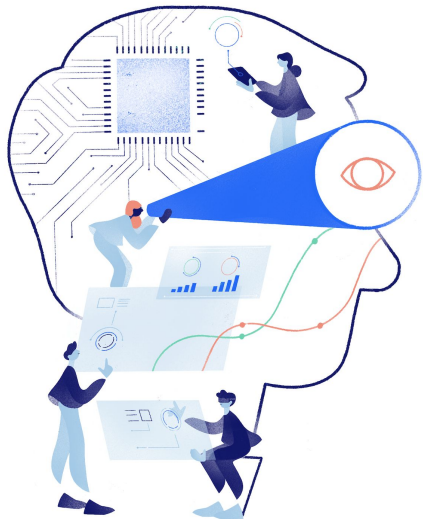- Training
- Coaching
- Team extensions

## AI strategy

- Discover the best opportunities in AI
- Value & effort estimation of AI opportunities

## Research

- Develop new AI models
- Explore new applications
- Validate research results

# Radix

## Natural Language Processing 🧠

- Text generation
- Question answering
- Summarization
- Named entity recognition
- Classification
- Duplicate detection
- Information retrieval

## Tabular Data 📊

- Classification
- Regression
- Anomaly detection
- Explainable prediction
- Confidence intervals
- Matrix imputation
- Feature engineering
- Data visualization

## Computer Vision 👁️

- Image detection
- Object detection
- Semantic segmentation
- Instance segmentation
- Real-time detection
- Hyperspectral imaging
- Super-resolution
- Optical Character Recognition

## Planning 📅

- Scheduling optimization
- Assignment optimization
- Inventory optimization
- Routing optimization

## Time Series 📈

- Forecasting
- Anomaly detection
- Classification
- Confidence intervals

## Speech 🎤

- Transcription
- Speech recognition
- Voice cloning
- Speaker identification

# Outline

- Who are InvestSuite & Radix
- **What are Large Language Models**
- The magic of LLMs
- Dispelling the magic of LLMs
- Applications of LLMs
- What does the future hold

# What is an LLM?

*The best thing about AI is its ability to*

| | |
|---|---|
| learn | 4.5% |
| predict | 3.5% |
| make | 3.2% |
| understand | 3.1% |
| do | 2.9% |

# What is an LLM?

The best thing about AI is its ability to create

The best thing about AI is its ability to create worlds

The best thing about AI is its ability to create worlds that

The best thing about AI is its ability to create worlds that are

The best thing about AI is its ability to create worlds that are both

The best thing about AI is its ability to create worlds that are both exciting

# Base LLM:

## Completes documents

Alice has 3 brothers. Each brother has 2 sisters. How many sisters does Alice have?

(A) 4

(B) 5

(C) 3

(D) 6

Submit    274

# Instruction tuned LLM:

## Answers questions

Alice has 3 brothers. Each brother has 2 sisters. How many sisters does Alice have?

Alice has 5 sisters.

4 = 75.00%

5 = 11.13%

2 = 9.64%

6 = 2.49%

four = 0.74%

Total: -2.20 logprob on 1 tokens
(99.00% probability covered in top 5 logits)

Submit          25

# ChatGPT:

## A web app on top of GPT-4

**You**
Alice has 3 brothers. Each brother has 2 sisters. How many sisters does Alice have?

**ChatGPT**
Alice has 2 sisters.

Message ChatGPT…

ChatGPT can make mistakes. Consider checking important information.

[1] ChatGPT: https://chat.openai.com

# Outline

Demo

**InvestSuite**

Investment
Platforms

Portfolio
Construction

Investment
Insights & Reporting

**Robo Advisor**  **Self Investor**   **Portfolio Optimizer**   **StoryTeller**  **Insight API**

**You**

| top_10_stocks_closing_prices.csv |
| Spreadsheet |

Here are some stock prices. What are the annualised return and risk of these stocks? Show me with a graph.

**ChatGPT**

⬤

Message ChatGPT...

ChatGPT can make mistakes. Consider checking important information.

Laurent Sorber

# The keys to success

**Instruction tuning**
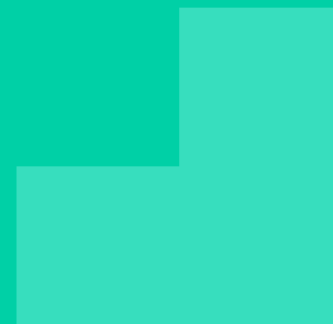
**Scaling to 1T tokens** (≈ words)

# Outline

- Who are InvestSuite & Radix
- What are Large Language Models
- The magic of LLMs
- **Dispelling the magic of LLMs**
- Applications of LLMs
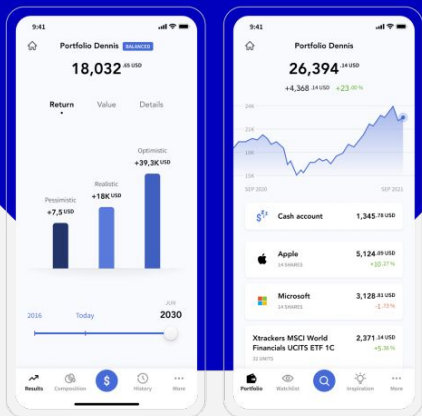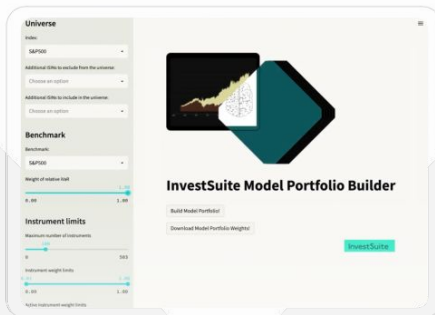- What does the future hold

# All LLMs are Transformers

| Input text | Tokenize | LLM | Output probabilities | Sample |
|---|---|---|---|---|

**Weights & Biases is the**

[1135, 2337, 1222, 8436, 1386, 318, 262]

$$\begin{bmatrix} 1 & 2 & 1 \\ 0 & 1 & 0 \\ 2 & 3 & 4 \end{bmatrix} \times \begin{bmatrix} 2 & 5 \\ 6 & 7 \\ 1 & 8 \end{bmatrix}$$

| ... | ... |
|---|---|
| We | 0.001 |
| ights | 0.001 |
| & | 0.001 |
| the | 0.001 |
| machine | 0.15 |
| learning | 0.001 |
| platform | 0.2 |

**machine**

[1] LLM image source: https://wandb.ai/darek/llmapps/reports/A-Gentle-Introduction-to-LLM-APIs--Vmlldzo0NjM0MTMz
[2] GPU image source: https://www.nvidia.com/en-us/data-center/h100/
[3] Matrix multiplication source: https://math.stackexchange.com/questions/3715598/when-should-hadamard-matrix-multiplication-be-used-vs-other-methods

# What do LLMs imitate?

✅ Cognitive ability          ❌ Consciousness

❌ Sentience          ❌ Will

# The T-shape of an LLM

Breadth of experience,
knowledge, and skills

- Hallucination
- Weak reasoning
- Weak planning

Depth of ability
in a given domain

LLM

Human expert

# The T-shape of an LLM

Breadth of experience, knowledge, and skills

✅ No hallucination

✅ Good reasoning

✅ Good planning

Depth of ability in a given domain

AGI ≥ median human depth

LLM

Human expert

# The T-shape of an LLM

Breadth of experience, knowledge, and skills

✅ No hallucination

✅ Exceptional reasoning

✅ Exceptional planning

Depth of ability in a given domain

**ASI** ≥ expert human depth

**LLM**

**Human expert**

# Outline

- Who are InvestSuite & Radix
- What are Large Language Models
- The magic of LLMs
- Dispelling the magic of LLMs
- **Applications of LLMs**
- What does the future hold

# The T-shape of an LLM app

Breadth of experience, knowledge, and skills

✅ No hallucination

✅ Guided reasoning

✅ Guided planning

Depth of ability in a given domain

☐ **LLM app** = median human depth

■ **LLM**

■ **Human expert**

# Anatomy of an LLM app

**LLM application**

**Software tools**
*Python, applications*

**Files**
*Confluence, Sharepoint*

**LLM**

Context window

**Audio & Video**
*Speech, images, video*

**Network**
*Databases, internet*

# LLM app: Portfolio Assistant

**LLM application**

Half the errors    Half the errors    Half the errors

100%

GPT-4

Task accuracy

+10%    +5%

GPT-3.5

0%

Base LLM | + Instruction tuning (SFT, RLHF) | + Zero-shot prompting[1] | + Few-shot prompting[2] | + Multi-step prompting[3,4,5] | + Retrieval | + Parameter Efficient Fine-tuning

★ GPT-4 achieves a few-shot chain-of-thought performance on MMLU of 86%: https://arxiv.org/abs/2305.17306
★ GPT-3.5 achieves a few-shot chain-of-thought performance on MMLU of 67%: https://arxiv.org/abs/2305.17306
★ GPT-3 achieves a few-shot performance on MMLU of 44%: https://arxiv.org/pdf/2009.03300.pdf
[1] An automatically discovered chain-of-thought prompt generalizes to novel models and datasets: https://arxiv.org/pdf/2305.02897.pdf
[2] Language Models are Few-Shot Learners: https://arxiv.org/pdf/2005.14165.pdf
[3,4,5] Self-consistency: https://arxiv.org/pdf/2203.11171.pdf, Self-reflection: https://evjang.com/2023/03/26/self-reflection.html, Tree of Thoughts: https://arxiv.org/pdf/2305.10601.pdf

* GPT-4 achieves a few-shot chain-of-thought performance on MMLU of 86%: https://arxiv.org/abs/2305.17306
* GPT-3.5 achieves a few-shot chain-of-thought performance on MMLU of 67%: https://arxiv.org/abs/2305.17306
* GPT-3 achieves a few-shot performance on MMLU of 44%: https://arxiv.org/pdf/2009.03300.pdf
[1] An automatically discovered chain-of-thought prompt generalizes to novel models and datasets: https://arxiv.org/pdf/2305.02897.pdf
[2] Language Models are Few-Shot Learners: https://arxiv.org/pdf/2005.14165.pdf
[3,4,5] Self-consistency: https://arxiv.org/pdf/2203.11171.pdf, Self-reflection: https://eyiang.com/2023/03/26/self-reflection.html, Tree of Thoughts: https://arxiv.org/pdf/2305.10601.pdf

# LLM app: StoryTeller podcast

**LLM application**

# LLM app: StoryTeller podcast

**LLM application**

Software tools
*Python, applications*

Files
*Confluence, Sharepoint*

**LLM**

Context window

**Audio & Video**
*Speech, images, video*

**Network**
*Databases, internet*

Spotify®

Apple Podcasts

```
---Start Podcast Script---

Adam: Welcome back to Planet Money. I'm Adam
Davidson, and with me as always, is Alex Blumberg.
How are you doing today, Alex?

Alex: Doing great, Adam. Ready to unpack some
finance stories for our listeners.

Adam: Excellent. So the mystery box we're opening
today is the ING Sustainable Balanced EUR R Cap
fund. Ready to dive in?

(Ambient instrumental- "intriguing" music playing)
```

# LLM app: Lexicon Assistant

**LLM application**

# LLM app: Lexicon Assistant

## LLM application

**Software tools**
*Python, applications*

**LLM**

Context window

**Files**
*Confluence, Sharepoint*

RAG

**Audio & Video**
*Speech, images, video*

**Network**
*Databases, internet*

PDF  PDF  PDF

---

**Understanding Volatility in Investments**

**Lexicon Assistant** ⌄

**You**
What does volatility mean?

**Lexicon Assistant**
Volatility is a statistical measure of the fluctuations in a security's price. It can also refer to the fluctuations in a particular market. When we say an investment or market has high volatility, it means the price of that investment or the values in that market can change a lot in a short period of time. This indicates a higher level of risk because the price can go up or down quite rapidly.
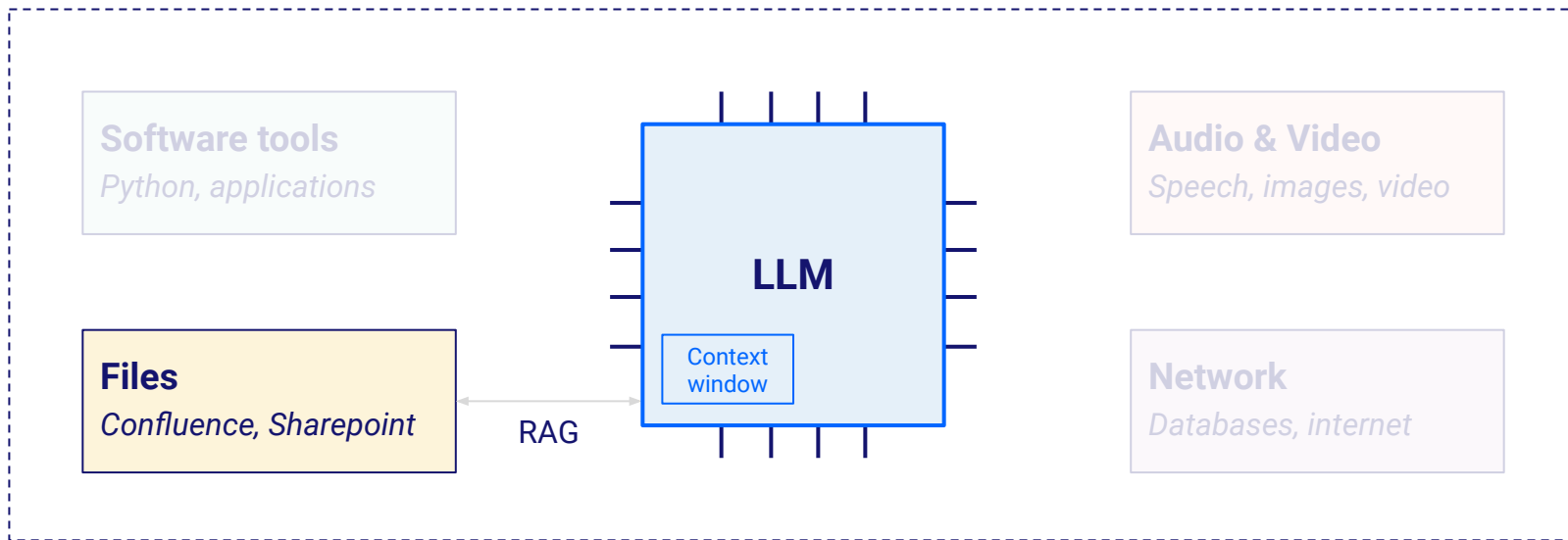
For example, think of a roller coaster – it goes up and down frequently and rapidly, which is similar to how a highly volatile investment behaves. On the other hand, a low volatility investment is more like a gentle car ride on a smooth road, with fewer ups and downs.

Message Lexicon Assistant...

ChatGPT can make mistakes. Consider checking important information.
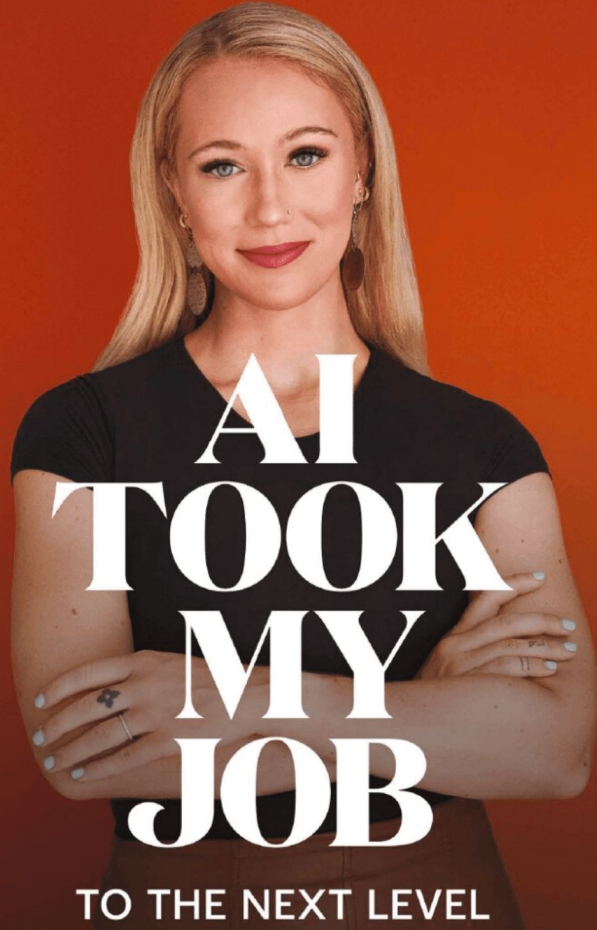
# Outline

- Who are InvestSuite & Radix
- What are Large Language Models
- The magic of LLMs
- Dispelling the magic of LLMs
- Applications of LLMs
- **What does the future hold**

Personal Computer

Talk to Pi,
your personal AI

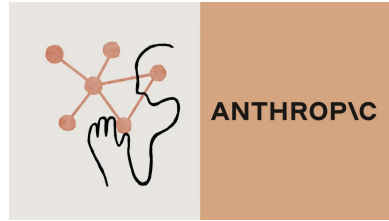heypi.com

Personal Intelligence

# AI political compass

# What are the H100s working on?



Hallucination

Alignment, Interpretability

Reasoning & planning

Reasoning & planning

Multi-modality

Robotics

Retrieval

Reasoning & planning

'Truth'

# AGI, when?



"Less than 10 years"[1]

Sam Altman (OpenAI)

"Within 18 months"[2]

Mustafa Suleyman (Inflection AI)

"As early as 2025, expected 2028"[3]

Shane Legg (Google DeepMind)

[1] Sam Altman (OpenAI) interview, 21 October 2023: https://www.youtube.com/watch?v=byYlC2cagLw
[2] Mustafa Suleyman (Inflection AI) interview, 14 November 2023: https://www.youtube.com/watch?v=VUOPCIgj-7U
[3] Shane Legg (Google DeepMind) interview, 26 October 2023: https://www.youtube.com/watch?v=Kc1atfJkiJU

# AGI, when?



"Less than 10 years"[1]

Sam Altman (OpenAI)

"Within 18 months"[2]

Mustafa Suleyman (Inflection AI)

"As early as 2025, expected 2028"[3]

Shane Legg (Google DeepMind)

[1] Sam Altman (OpenAI) interview, 21 October 2023: https://www.youtube.com/watch?v=byYlC2cagLw
[2] Mustafa Suleyman (Inflection AI) interview, 14 November 2023: https://www.youtube.com/watch?v=VUOPCIgj-7U
[3] Shane Legg (Google DeepMind) interview, 26 October 2023: https://www.youtube.com/watch?v=Kc1atfJkiJU

# Thank you!

InvestSuite

Radix