

# Decoding Al series

### Session 2

Unleashing the power of Large Language Models

Part 2 – Workflows and complex interactions

### Our awesome speaker



### **Gil Fernandes**

**AI Solutions Engineer** 



Welcome

## From Roots to Fruits

• Journey to create value from LLMs





Agenda



Insights from previous webinar	<u>04</u>
In-context Learning	<u>11</u>
RAG (Retrieval Augmented Generation)	<u>14</u>
Reasoning strategies	<u>18</u>
D-Wise Demo	<u>46</u>
Credits	<u>48</u>
Thank you for joining	<u>49</u>



parameter, and provide par-

# E Audience poll

## For which type of tasks are you using LLMs and MLLMs?

- 1. Write communication, like e.g. emails, memos
- 2. Improve and correct text
- 3. Translate text
- 4. Summarise text



- 5. Generate images
- 6. Generate audio files, like music
- 7. Generate video
- 8. Generate code (programming)

LLMs suited for many scenarios, not for others

LLMs suited for:

Great at a variety of NLP tasks

Converting unstructured data into structured data

Tackling the problem of "dark data"





www.onepointltd.com/techtalk





LLMs suited for many scenarios, not for others

LLMs struggle with:

Limited reasoning

Long term planning

Need to be "grounded" (hallucinations)



Knowledge cut-offs

Missing sources

Refusals



www.onepointltd.com/techtalk





## Getting value out of LLMs

Improve direct interaction with LLMs

Get better results with prompting

Leverage Dark Data

Create translation layers with simple agents



Working with unstructured data

Create automation with LLM driven workflows

Simulate complex scenarios with agent ensembles



www.onepointltd.com/techtalk











# How we interact with LLMs (possibilities)



![](_page_9_Picture_5.jpeg)

![](_page_9_Picture_6.jpeg)

# In-context learning

![](_page_10_Picture_1.jpeg)

A LINE AND AND A DESCRIPTION OF A DESCRIPT

preservation president presidents

A REPORT OF A DESCRIPTION OF A DESCRIPTI

# LLMs can learn on the fly

LLMs can learn from a given text and answer questions about it.

Based on the following text can you describe the Mamba neural network architecture?

•••

Foundation models (FMs), or large models pretrained on massive data then adapted for downstream tasks, have emerged as an effective paradigm in modern machine learning. The backbone of these FMs are often sequence models, operating on arbitrary sequences of inputs from a wide variety of domains such as language, images, speech, audio, time series, and genomics \parencitesutskever2014sequence,dosovitskiy2020image,oord2016wa venet,brown2020language,ismail2019deep,poli2023hyena. While this concept is agnostic to a particular choice of model architecture, modern FMs are predominantly based on a single type of sequence model: the Transformer \parencitevaswani2017attention and its core attention layer \parencitebahdanau2015neural The efficacy of self-

![](_page_11_Picture_6.jpeg)

The Mamba neural network architecture represents a significant advancement in sequence modeling, designed to combine the strengths of structured state space models (SSMs) with the capabilities of modern Transformers, while addressing some of their key limitations. Here's a detailed description of its architecture and features:

### Key Components of Mamba Neural Network Architecture

- 1. Selective State Space Models (SSMs):
  - Selection Mechanism: Mamba incorporates a novel selection mechanism that allows the model to filter and prioritize relevant information based on the input. This mechanism is parameterized by the input, enabling the model to ignore irrelevant data and retain important information indefinitely.
  - Hardware-aware Algorithm: To maintain computational efficiency, Mamba uses a hardware-aware algorithm that processes the model recurrently with a scan operation instead of convolution. This avoids unnecessary IO access between different levels of GPU memory, ensuring faster computations. The model scales linearly with sequence length and performs efficiently on modern hardware like A100 GPUs.

Source: ChatGPT-4o

![](_page_11_Picture_13.jpeg)

In-context learning

## LLMs can learn on the fly

Why is this extremely useful?

Allows LLMs to access private knowledge bases or very recent knowledge

Allows the creation of Retrieval Augmented Generation (RAG) applications

![](_page_12_Picture_5.jpeg)

**Reduces hallucinations** 

Allows to trace back sources

![](_page_12_Picture_8.jpeg)

# RAG (Retrieval Augmented Generation)

![](_page_13_Picture_1.jpeg)

A Constant ( I as a reason of B Constant for the second seco

O DE CONTRA DE CONTRA

general space, produced

RAG

## **Retrieval Augmented Generation**

## Enhance LLM knowledge with private data and vector databases

Involved components

LLM

Vector database

Orchestration

User sends request

![](_page_14_Picture_8.jpeg)

![](_page_14_Figure_9.jpeg)

![](_page_14_Figure_10.jpeg)

![](_page_14_Picture_11.jpeg)

RAG

## **Retrieval Augmented Generation**

RAG data needs to be "indexed", i.e. converted to vectors which have to be generated by the LLM.

![](_page_15_Figure_3.jpeg)

![](_page_15_Picture_4.jpeg)

![](_page_15_Picture_5.jpeg)

## RAG Business use cases

Talk to your documents

Product search

**Content creation** 

![](_page_16_Picture_5.jpeg)

![](_page_16_Figure_6.jpeg)

![](_page_16_Picture_7.jpeg)

Reasoning strategies

![](_page_17_Picture_1.jpeg)

Chain of Thoughts

Tree of Thoughts

**Buffer of Thoughts** 

Monte Carlo Tree Self-refine

Onepoint TechTalk

A LINE AND AND A DESCRIPTION OF A DESCRIPT

preservation president presidents

A REPORT OF A DESCRIPTION OF A DESCRIPTI

## Step-by-step reasoning

The main idea is to formulate a problem and instruct the model to break it down in sub-problems and solve these until the final answer is found

The prompt will have two parts:

Describe the main problem

![](_page_19_Picture_5.jpeg)

Instructions on breaking down and solving each sub-problem until the main problem is solved

![](_page_19_Picture_7.jpeg)

![](_page_20_Figure_1.jpeg)

![](_page_20_Figure_2.jpeg)

![](_page_20_Picture_3.jpeg)

![](_page_20_Picture_4.jpeg)

## Example – Step 1

![](_page_21_Picture_2.jpeg)

- The song you're referring to is "Come On Eileen" by Dexys Midnight Runners. \$
  - type of medication).
  - Song: "Come On Eileen".
  - Female name: Eileen.
  - adverb, but it carries an urgent, sometimes inconsistent connotation.
  - somewhat inconsistent in its plea.

So, the song is "Come On Eileen" by Dexys Midnight Runners.

Source: ChatGPT-4o

![](_page_21_Picture_12.jpeg)

James likes to hear music from the Eighties and there is one band which he appreciates very much. They have this song in which the name of a female is related to an adverb. The pronoun is related to inconsistency and the person in the song is female. The band's name is somehow related to healing. Which song are we talking about?

• Band: Dexys Midnight Runners (related to healing as "Dexys" refers to dextroamphetamine, a

• Adverb: The word "Come" in "Come On Eileen" can be seen as an imperative verb, not an

• Pronoun: "Come on" implies a direct address to Eileen, and the call to action can be seen as

## Example – Step 2

![](_page_22_Picture_2.jpeg)

the fact that the name is related to healing? Given the period (the Eighties) and the clue that the band's name is related to healing, the band in question is likely The Cure.

• Band: The Cure.

\$

• Song: "Why Can't I Be You?"

In this song, the name of the person isn't explicitly stated, but the pronoun "you" relates to inconsistency (the longing to be someone else), and the song's title includes an interrogative adverb ("Why"). However, it seems to miss a specific female name and the adverb connection.

Revisiting your criteria more thoroughly:

- Band: The Cure.
- Song: "Close to Me" from The Cure.

Although "Close to Me" does not explicitly mention a female name related to an adverb, the pronoun "me" is used, and the song revolves around personal feelings, which can be inconsistent.

Source: ChatGPT-40

![](_page_22_Picture_12.jpeg)

Can you first try to identify the band by looking at the period and then

Now, considering the song by The Cure that fits your description:

![](_page_22_Picture_17.jpeg)

## Example – Step 3

![](_page_23_Picture_2.jpeg)

\$

Given that the band is The Cure and we are looking for a song from the Eighties that features a female name related to an adverb and a pronoun implying inconsistency, the song in question is likely "Charlotte Sometimes."

- Band: The Cure
- Song: "Charlotte Sometimes"
- Female name: Charlotte
- Adverb: Sometimes (related to inconsistency)

The song "Charlotte Sometimes" by The Cure fits all the criteria perfectly.

Source: ChatGPT-4o

![](_page_23_Picture_12.jpeg)

You got that right. It is the cure. Now focus on the name of the song. There is a female name in the song. Which could it be?

• **Pronoun**: The song's narrative implies varying perspectives and feelings.

## Induced chain of thoughts

What if we could let the model create on its own the chain of thoughts?

You can use a prompt to induce this behaviour using **custom instructions**.

![](_page_24_Picture_4.jpeg)

## Induced chain of thoughts

Example:

You carefully provide factual, accurate, thoughtful and nuanced answers. You are brilliant at reasoning. If you think that you might not be correct, you say so.

You always spend a few sentences explaining the background context, assumptions and step-by-step thinking before you try to answer a question.

![](_page_25_Picture_5.jpeg)

Mary has 3 brothers. Each of her brothers has 3 sisters. How many sisters does Mary have?

B

Let's break down the information given in the problem:

- 1. Mary has 3 brothers.
- 2. Each of her brothers has 3 sisters.

Given that Mary is one of those sisters, the brothers share the same sisters. Therefore, Mary and her brothers share the same siblings.

Since each brother has 3 sisters and Mary is one of those sisters, it means Mary has 2 other sisters.

So, Mary has 2 sisters.

Source: ChatGPT-4o

![](_page_25_Picture_15.jpeg)

## Induced chain of thoughts

Example:

You carefully provide factual, accurate, thoughtful and nuanced answers. You are brilliant at reasoning. If you think that you might not be correct, you say so.

You always spend a few sentences explaining the background context, assumptions and step-by-step thinking before you try to answer a question.

![](_page_26_Picture_5.jpeg)

Mary has 3 brothers. Each of her brothers has 3 sisters. How many sisters does Mary have?

B

Let's break down the information given in the problem:

- 1. Mary has 3 brothers.
- 2. Each of her brothers has 3 sisters.

Given that Mary is one of those sisters, the brothers share the same sisters. Therefore, Mary and her brothers share the same siblings.

Since each brother has 3 sisters and Mary is one of those sisters, it means Mary has 2 other sisters.

So, Mary has 2 sisters.

Source: ChatGPT-4o

![](_page_26_Picture_15.jpeg)

# ChatGPT o1-preview and o1-mini

Operate in 2 steps: reasoning and response.

Reasoning step creates action plan

Response step generates the answer

![](_page_27_Picture_5.jpeg)

![](_page_27_Picture_6.jpeg)

## ChatGPT o1-preview and o1-mini

Experience: slower, but much better at reasoning

However it still makes mistakes, so be aware of that.

![](_page_28_Picture_4.jpeg)

### Back to Game of 24. Can you solve it with these numbers: 5, 3, 5, 3?

ChatGPT o1-mini ~

#### **Alternative Representation:**

Another way to represent the same solution is:

$$(3 imes 5) imes \left(rac{5+3}{5}
ight)=15 imes rac{8}{5}=15 imes 1.6=24$$

#### **Conclusion:**

By strategically combining addition, division, and multiplication, we've successfully reached the target number **24** using the numbers **5**, **3**, **5**, **3**. This solution adheres to the rules of the Game of 24 by using each number exactly once and employing standard arithmetic operations.

#### Answer:

Yes, one way to make 24 using 5, 3, 5, and 3 is:  $3 imes5 imes\left(rac{5+3}{5}
ight)=24$ 

() በ በ የ 2 ~

Source: ChatGPT o1-mini

![](_page_28_Picture_16.jpeg)

![](_page_29_Picture_1.jpeg)

A LINE AND AND A DESCRIPTION OF A DESCRIPT

preservation president presidents

A REPORT OF A DESCRIPTION OF A DESCRIPTI

![](_page_29_Picture_12.jpeg)

## Step-by-step with multiple options

### Instruct the LLM to break down the problem in steps and then to solve each step iteratively.

![](_page_30_Figure_3.jpeg)

Source: <u>https://github.com/princeton-nlp/tree-of-thought-llm</u>

![](_page_30_Picture_5.jpeg)

![](_page_30_Picture_6.jpeg)

# E Audience poll

Please detect the game of 24 combination that cannot be solved For example for 6, 2, 16, 6 can be solved. Possible solution is (6 \* 16) / (6 - 2)

- 1. 7, 2, 16, -6
- 2. 5, 3, 3, 6
- 3. 6, 6, 6, 2

![](_page_31_Picture_5.jpeg)

- 4. 5, 5, 5, 7
- 5. All are right
- 6. Don't know...l give up

![](_page_31_Picture_10.jpeg)

## Game of 24 – simple prompts go wrong (1 of 2)

![](_page_32_Picture_2.jpeg)

Source: <u>https://www.youtube.com/watch?v=fWEEu9JxyYo&t=28s</u> How to play a math game

![](_page_32_Picture_4.jpeg)

![](_page_32_Figure_5.jpeg)

Source: ChatGPT-4o

![](_page_32_Picture_7.jpeg)

## Game of 24 – simple prompts go wrong (2 of 2)

![](_page_33_Picture_2.jpeg)

Source: <a href="https://www.youtube.com/watch?v=fWEEu9JxyYo&t=28s">https://www.youtube.com/watch?v=fWEEu9JxyYo&t=28s</a> How to play a math game

![](_page_33_Picture_4.jpeg)

And a solution for 7, 2, -6, 16?

For the numbers **7**, **2**, **-6**, **and 16**, let's work through possible operations to reach 24.

One possible solution is:

$$16 - (7 - 2) - (-6) = 16 - 5 + 6 = 24$$

### Step-by-step breakdown:

- 1. Subtract 7 from 2 to get 7-2=5.
- 2. Subtract this result from 16: 16 5 = 11.
- 3. Subtract -6 from the result (which is the same as adding 6): 11+6=24.

Thus, this combination satisfies the game's rules and results in 24.

Source: ChatGPT-4o

\$

![](_page_33_Figure_15.jpeg)

## The idea behind Tree of Thoughts

Tree of Thoughts is an iterative process which allows to solve problems in a combinatorial problem space.

At each step one number disappears from the equation until we have a solution.

![](_page_34_Picture_4.jpeg)

![](_page_34_Figure_6.jpeg)

Source: https://github.com/princeton-nlp/tree-of-thought-llm

## The idea behind Tree of Thoughts

### Propose

Create potential solutions for the current step

### Values

Calculate the values for the potential solutions

![](_page_35_Picture_6.jpeg)

### Evaluate

Evaluate the success probability for each solution

Sample

Get the best evaluated results

![](_page_35_Picture_12.jpeg)

## The workflow

![](_page_36_Figure_2.jpeg)

Source: https://github.com/princeton-nlp/tree-of-thought-llm

![](_page_36_Picture_4.jpeg)

![](_page_36_Picture_5.jpeg)

## Performance ToT versus other strategies

Method	Success	(a) S
IO prompt CoT prompt CoT-SC (k=100) ToT (ours) (b=1) ToT (ours) (b=5)	7.3% 4.0% 9.0% 45% <b>74%</b>	0.6
IO + Refine (k=10) IO (best of 100) CoT (best of 100)	27% 33% 49%	0.2

Source: https://github.com/princeton-nlp/tree-of-thought-llm

![](_page_37_Picture_4.jpeg)

![](_page_37_Figure_5.jpeg)

38

# Buffer of Thoughts

![](_page_38_Picture_1.jpeg)

A LINE AND AND A DESCRIPTION OF A DESCRIPT

preservation president presidents

A REPORT OF A DESCRIPTION OF A DESCRIPTI

Buffer of Thoughts

## Performance ToT versus other strategies

![](_page_39_Figure_2.jpeg)

![](_page_39_Picture_5.jpeg)

![](_page_39_Picture_7.jpeg)

Buffer of Thoughts

## Detailed workflow

Distill problem parameters

Lookup template (prompt template) in buffer

Combine template and parameters

Execute template

![](_page_40_Picture_6.jpeg)

![](_page_40_Picture_7.jpeg)

![](_page_40_Figure_8.jpeg)

![](_page_40_Picture_9.jpeg)

# Monte Carlo Tree Self-refine

THE R. P. LEWIS CO., LANSING MICH.

![](_page_41_Picture_2.jpeg)

STREPENSER'S CONTRACTOR STREPENSES STREPENSES STREPENSES STREPENSES

A Developed ( ) has present and the B Developed Pro-

Concession of the Print of the Number of Street, or other Designation of the Print of the Print

NAME AND ADDRESS OF TAXABLE PARTY.

A REAL PROPERTY AND ADDRESS OF THE OWNER, THE PARTY OF TH

parameter at a paral parameter para

Monte Carlo Tree Self-refine

Techniques for improved reasoning with LLMs (1 of 2)

Boosts the reasoning capabilities of LLMs

Help to solve very complex problems with even "small" LLMs.

Using a recursive answer improvement and rating pipeline you can create reasoning logical "trees" which boost the model's reasoning capabilities.

![](_page_42_Picture_5.jpeg)

![](_page_42_Figure_7.jpeg)

Source: www.youtube.com/watch?v=mfAV bigdRA

![](_page_42_Picture_9.jpeg)

### Monte Carlo Tree Self-refine

## Techniques for improved reasoning with LLMs (2 of 2)

![](_page_43_Figure_2.jpeg)

Source: arxiv.org/abs/2406.07394

![](_page_43_Picture_4.jpeg)

![](_page_44_Picture_0.jpeg)

7, 2, 16, -6 (2\*7) + (16-6)

5, 3, 3, 6 ((5 \* 3) + 3) + 6

6, 6, 6, 2

(6 \* 6) - (6 \* 2)

× 5, 5, 5, 7

Does not form 24

![](_page_44_Picture_7.jpeg)

![](_page_44_Figure_8.jpeg)

![](_page_44_Figure_9.jpeg)

# Live demo

![](_page_45_Picture_1.jpeg)

A LINE AND AND A DESCRIPTION OF A DESCRIPT

preservation president presidents

A REPORT OF A DESCRIPTION OF A DESCRIPTI

Free AI tool

# Explore your data health with Onepoint D-Wise<sup>™</sup>

![](_page_46_Picture_2.jpeg)

Embark on your data wellness journey with our free tool, powered by generative AI and the Onepoint Data & Analytics Body of Knowledge™

![](_page_46_Picture_4.jpeg)

![](_page_46_Picture_5.jpeg)

![](_page_46_Figure_6.jpeg)

Scan the QR code or visit <u>www.onepointltd.com/data-wellness</u>

![](_page_46_Picture_8.jpeg)

## Credits

![](_page_47_Picture_1.jpeg)

### Tree of Thoughts

chatgpt.com

github.com/princeton-nlp/tree-of-thought-llm

arxiv.org/pdf/2406.04271

## Accessing GPT-4 level Mathematical Olympiad Solutions via Monte Carlo Tree Self-refine with LLaMa-38B

arxiv.org/pdf/2406.07394 & www.youtube.com/watch?v=mfAV\_bigdRA

![](_page_47_Picture_9.jpeg)

### **Buffer of Thoughts**

48

## Thank you for joining

Please feel free to contact Gil Fernandes if you have any feedback about the session.

![](_page_48_Picture_2.jpeg)

Email at <u>techtalk@onepointltd.com</u>

Connect on LinkedIn <u>www.linkedin.com/in/gil-palma-fernandes</u>

Find Gil's Reflections on AI at <u>medium.com/@gil.fernandes</u>

![](_page_48_Picture_6.jpeg)

![](_page_48_Picture_7.jpeg)

![](_page_49_Picture_0.jpeg)

# Next webinar coming soon!

Details will be announced soon. Stay tuned!

Onepoint is a **boutique**, values-driven, businessoriented technology consultancy.

We architect, prototype, build, and manage data and AI powered solutions. We partner with global clients looking for high-impact, enterprise-grade advice and IT services to realise their most critical digital transformations.

London | Manchester | Pune

## Onepoint

Your trusted companions for the digital journey<sup>m</sup>

![](_page_50_Picture_5.jpeg)

onepointltd.com

hello@onepointItd.com

![](_page_50_Picture_8.jpeg)

![](_page_50_Picture_9.jpeg)

© Copyright 2024. Onepoint Consulting Ltd is a company registered in England. Company registration number 5516457. VAT registration number GB866120039. ISO27001 certified. Certified in the UK as an ethnic minority-owned business.