# Initial Steps in Integrating Large Reasoning and Action Models for Service Composition

Ilche Georgievski and Marco Aiello, IAAS-SC, University of Stuttgart, DE

SummerSOC 2025

# AI is eating the world

| | | |
|---|---|---|
| OpenAI | Cohere | Mistral AI |
| Google | Microsoft | Anthropic |
| DeepSeek | Meta | 01.AI |
| Inflection AI | Nvidia | Contextual AI, Inc. |
| Databricks | GPT-4 | Mosaic ML, Inc. |
| AWS | Claude | IBM |
| Meta AI | Why Labs | Alibaba |
| Amazon | Claude 3 Haiku | Falcon |

# AI is eating the world

**OpenAI**

estimated at 300B

| Rank | | Name | Market Cap | Price | Today | Price (30 days) | Country |
|------|---|------|------------|-------|-------|-----------------|---------|
| ⌃1 | 1 | Microsoft<br>MSFT | $3.476 T | $467.68 | ▲ 0.82% | | 🇺🇸 USA |
| ⌄1 | 2 | NVIDIA<br>NVDA | $3.414 T | $139.99 | ▼ 1.36% | | 🇺🇸 USA |
| | 3 | Apple<br>AAPL | $2.996 T | $200.63 | ▼ 1.08% | | 🇺🇸 USA |
| | 4 | Alphabet (Google)<br>GOOG | $2.049 T | $169.81 | ▲ 0.25% | | 🇺🇸 USA |
| | 5 | Meta Platforms (Facebook)<br>META | $1.721 T | $684.62 | ▼ 0.48% | | 🇺🇸 USA |
| | 6 | Tesla<br>TSLA | $917.00 B | $284.70 | ▼ 14.26% | | 🇺🇸 USA |
| | 7 | Oracle<br>ORCL | $479.91 B | $171.14 | ▲ 1.81% | | 🇺🇸 USA |
| | 8 | Palantir<br>PLTR | $282.97 B | $119.91 | ▼ 7.77% | | 🇺🇸 USA |
| | 9 | IBM<br>IBM | $248.01 B | $266.86 | ▲ 0.50% | | 🇺🇸 USA |
| | 10 | Adobe<br>ADBE | $176.95 B | $415.20 | ▲ 0.31% | | 🇺🇸 USA |
| | 11 | CoreWeave<br>CRWV | $64.82 B | $135.05 | ▼ 17.20% | | 🇺🇸 USA |
| | 12 | Cambricon Technologies<br>688256.SS | $35.43 B | $84.88 | ▼ 1.59% | | 🇨🇳 China |
| | 13 | Dynatrace<br>DT | $16.46 B | $54.94 | ▲ 1.55% | | 🇺🇸 USA |
| | 14 | Mobileye<br>MBLY | $13.39 B | $16.50 | ▼ 1.70% | | 🇮🇱 Israel |
| | 15 | Tempus AI<br>TEM | $10.15 B | $58.66 | ▼ 6.56% | | 🇺🇸 USA |
| | 16 | Aurora Innovation<br>AUR | $9.99 B | $5.65 | ▼ 2.75% | | 🇺🇸 USA |
| | 17 | UiPath<br>PATH | $7.09 B | $13.26 | ▲ 1.61% | | 🇺🇸 USA |

# "Company with no business plan buys company with no product"
## For 6.5 billion dollars







Jony Ive leading OpenAI's design work following $6.5B acquisition

TechCrunch

# Builder.ai

## AI for software engineering or 700 SE humans



promised to make software creation "as easy as ordering pizza"

raised $450 million and achieved a valuation of $1.5 billion

reportedly owes $85 million to Amazon and $30 million to Microsoft in unpaid cloud services

# on the other hand

much more sustainable than real AI...

To estimate the **daily electricity consumption** for 700 people:

- **Annual per capita consumption**: 1,395 kWh

- **Daily per capita consumption**:
  1,395 kWh ÷ 365 days ≈ **3.82 kWh/day**

- **Total for 700 people**:
  3.82 kWh/day × 700 people ≈ **2,674 kWh/day**

# Arxiv Papers having LLM or GPT in the title (all disciplines)

Last 12 months
17844

Last 12 months
17,844

18844

14133

9422

4711

0

2        8        7        25       64       184      3017     11295

2017    2018    2019    2020    2021    2022    2023    2024

©Marco Aiello, 2025

*based on arxiv publications*

# The AI Revolution and the Future of Work: Threats and Opportunities

**Marco Aiello**, University of Stuttgart

*AI is transforming the global labor market, signaling a shift where jobs are both displaced and newly created. But is it really creating enough new jobs?*

Technological advancement brings new entrepreneurial opportunities and has the power to revolutionize markets. In doing so, jobs can be swapped away, considerably downsized, or deeply changed. The neoclassical view on technological advancements postulates that whenever a certain type of job is made useless by technological innovation, new jobs appear on the market to compensate for the loss, possibly after a painful transition.[1] If the Second Agricultural Revolution finally resulted in workers moving from the fields to the factories, the end of the First Industrial Revolution marked the shift of work from being available in factories to being in offices. Usually, the number of office jobs created was greater than those lost. These revolutions brought economic growth in terms of average income per person and gross domestic product (GDP). Up until recently, economic advancement has guaranteed that labor markets have many new opportunities and that unemployment rates are low during periods of growth. Now, though, things appear to be different, as pointed out by several economists (for example, Paul Krugman[2]) and computer scientists (see Moshe Vardi[1]). Such a view is often referred to...

---

# A Paradigm Shift in Service Research: The Case of Service Composition

Marco Aiello, *Senior Member, IEEE*

**Abstract**—Recent advancements in artificial intelligence, particularly in machine learning and neural networks, have significantly influenced various domains, including service computing. Large Language Models (LLMs) are at the forefront of this transformation, introducing new paradigms for automation and decision-making. This paper examines the evolving impact of LLMs on service composition, a fundamental problem in service computing. By analyzing shifts in research approaches, methodologies, and system architectures, we highlight how LLM-driven automation challenges traditional composition techniques. The discussion provides insights into emerging opportunities, limitations, and research directions, emphasizing the need to rethink service composition in the era of AI-driven automation.

**Index Terms**—Service Composition, LLM, Philosophy of Science

## I. INTRODUCTION

Are we in the midst of a paradigm shift in Service Computing and, more generally, in Computing? As one does these days, to answer this question, I turn to prominent LLMs available online and ask the question. ChatGPT has more arguments in favor of saying that we are; Gemini is unsure and states that "the jury is out," and finally, Perplexity answers with a definite yes. But when can we say that a paradigm shift has occurred?

The term dates back to the best-selling book on the philosophy of science by Kuhn [1]. In his view, paradigm shifts mark moments when established science practices are replaced by new models and frameworks that redefine an entire field. From analyzing the history of natural science, Kuhn posits that science is cyclic, one phase of which is characterized by rapid and disruptive shifts marking the beginnings of new eras. Preliminary to the shift phase is a crisis one, in which the discipline under transformation is characterized by animated debate and the recurring to philosophy to frame the disruptive novelty. Such debate appears to be very current in our field, where, on the one side, a number of researchers identify elements of general intelligence and the ability to reason about reasoning (i.e., possessing a theory of mind) in current LLM systems [2], [3]. On the other hand, others are skeptical of their actual learning abilities. Noam Chomsky, well-known for his pioneering work on grammars, objects that LLMs do not learn grammar and do not learn as a human child would [4]. Others

Marco Aiello is with the Institute of Architecture of Application Systems, University of Stuttgart, Germany. e-mail: aiello@iee.org. The paper is based on the keynote delivered at IEEE ICWS 2024, Shenzhen, China, a video of which is available at https://tinyurl.com/2nkm5exv. I am grateful to Frank Leymann and the anonymous reviewers for suggestions that improved the paper; any remaining errors or oversights are solely my own. Figure 2 is based on the input of Yuchen Liu, whom I thank.

point out that LLMs utter statements without having a model of truth, ultimately "bullshitting" [5]. Either way, the current advancements in AI have taken our field and that of scientific research by storm. In the last 12 months, the number of papers posted on arxiv.org that have the term LLM or GPT in the title has reached almost 10,000 units; and AI use in research has seen rapid tenfold growth in all fields [6].

## II. SERVICE COMPOSITION: THEN AND NOW

To make things more concrete, let me focus on service computing and the central problem of service composition, a field I have been active in for more than two decades. The idea of service composition is that functionalities of all kinds are available via standard interfaces on a network and can be used anytime on a per-need basis to create dynamic, scalable, QoS-aware, added-value systems and services [7] [8]. The first works on this problem appeared at the beginning of the century, with Artificial Intelligence Planning as one of the preferred tools to achieve it [9]. The AI planning-based solutions, while very elegant, are brittle. To function, the systems need semantic annotations and significant human effort to provide tailor-made knowledge in addition to the availability of the services. These types of solutions had a similar destiny to that of the Semantic Web. They were nice in theory, but there was too much handwork to be adopted. It should be noticed that the composition so made are correct by design, as one can prove that the composition is a realization of the formal specifications provided by the user, e.g., [10].

LLMs can also be used to create service compositions. One can express a goal in natural language and ask to produce orchestrations [11]. Initial results show great promise as the LLMs are capable of generating syntactically correct code, the logic of which represents the composition requested well. Though the output is not necessarily executable, requiring manual intervention, especially for hallucinated APIs. Here, we see the first differences in the current shift. Traditional approaches to service composition require a great deal of effort in designing the architecture, providing domain knowledge in a formal way, and being applicable to specific domains. In addition, the system could only be designed and used by highly specialized experts. In contrast, LLM-based approaches are generally applicable given their foundation models, require no encoding of domain knowledge, and can be broadly applied without modifications. Furthermore, non-experts can use them to a large extent. The major drawback is that the solutions are not necessarily correct and ready to run.

Looking more in detail at how service composition architectures are changing between traditional approaches and

---

# A Challenge for the Next 50 Years of Automated Service Composition

Marco Aiello[✉]

Department of Service Computing, IAAS, University of Stuttgart, Stuttgart, Germany
marco.aiello@iaas.uni-stuttgart.de

**Abstract.** Automated Service Composition emerged as a promising area of research at the beginning of the century. After twenty years, it appears to have reached a stagnating state where only little progress is made. In the present vision paper, I propose a challenge for automated service composition to be achieved in the next 50 years. I set a scene in 2052 that service composition should be able to handle by then. Finally, I draw a parallel with autonomous driving to identify the major milestones in the quest to fully autonomous service composition.

**Keywords:** Automated Service Composition · Service-oriented Computing · Maturity Levels · Artificial Intelligence Planning

## 1 The Promises of Automated Service Composition

Automated service composition refers to systems that utilize distributed, discrete units of software by orderly invoking their execution with the goal of satisfying a set of user-defined specifications. The core idea is as old as the field of software engineering. In fact, as soon as software was complex enough to require artisan talent and engineering techniques, the intuition of using modular designs came about. Instead of writing code for every subtask, one could reuse parts of existing code, possibly resident remotely on a network. To make things simpler for the developer, the input/output syntax of these parts must be precisely specified in order to enable composition. These were the first steps in the direction of manual software composition and, with the subsequent advent of software services as units of invokable functionalities, of service composition.

### 1.1 A Parallel with the Automotive Industry

Automated service composition refers to systems... [continued]

try and the process of driving a car will help my position evolution. Since the first 'ride' of have been manually steering their vehicles . With the passing of time, more automation driver, such as synchronised gear shifting,

---

**EDITORIAL**

# Service composition in the ChatGPT era

Marco Aiello[1] · Ilche Georgievski[1]

ChatGPT recently attracted vast attention in and outside the research community for its conversational abilities that mimic human ones exceptionally well. At the heart of systems like ChatGPT are Large Language Models (LLM). These models, rooted in deep neural networks, have the ability to predict the next textual token in a series of tokens based on statistical occurrences in extremely large data sets [1]. When the models are sufficiently big and well-tuned, one observes the "unreasonable effectiveness of data" [2] in how the system generates perfectly intelligible and believable sentences. Such ability to have human-like conversations with a software system is both stunning for the quality of the conversation and mind-blowing in terms of the potential impact on society and the job market in particular [3, 4]. There is controversy on whether or not such systems manifest forms of artificial intelligence. Researchers at Microsoft, for instance, attribute signs of intelligence to the current fourth version of Generative Pre-trained Transformer (GPT-4), which is in development at the time of writing [5]. Some authors have successfully solved Theory of Mind tasks using such tools. Kosinski reports a success rate of 95% using GPT-4 in solving false-belief tasks. Other authors are more careful with ... anthropomorphization of ChatGPT-like systems [6, 7]. What is sure is that the embedding of an LLM into a system makes it a very powerful tool. Of interest to us in this editorial is the LLM capability to generate programs [8] ... potential impact on Service-Oriented Computing and Applications.

The problem of automated service composition is central to the field of Service-Oriented Computing and Applications. It is the unsupervised, automated composition of services available on a network in order to execute any task relying on multiple, loosely coupled services, possibly without prior knowledge of such services. This would guarantee that virtually any task can be executed by relying on third-party implementations and resources.

Such a vision of automation is rooted in the fields of software engineering and component-based software engineering. It emerged in a fervent moment of technological evolution. At the beginning of the century, the Internet was becoming pervasive, and the Web emerged as a central technology also for businesses. Every company was moving toward having a web server ... enterprise information ... ation was also ... being stripped ... turned into a ... XML. Specifi ... wiring of m ... describing se ... trations (BPI ... also propose ... such technolo ... be discovered ... (*publish-find* ...

The techno ... inspired man ... ing automate ... survey of ... emerged from ... stand what a ... the signature ... the exchange ... mentation an ... composition ... tems were ju ... could not rea ... would do, or ... the services h ... and the sema ... that made th ... the road of ...

---

[✉] Marco Aiello
marco.aiello@iaas.uni-stuttgart.de

Ilche Georgievski
ilche.georgievski@iaas.uni-stuttgart.de

[1] Service Computing Department, IAAS, University of Stuttgart, Universitätsstr. 38, 70569 Stuttgart, BW, Germany

---

# Compositio Prompto: An Architecture to Employ Large Language Models in Automated Service Computing

Robin D. Pesl[1], Carolin Mombrey[1,3], Kevin Klein[1,2], Denesa Zyberaj[1,2], Ilche Georgievski[1], Steffen Becker[1], Georg Herzwurm[1], and Marco Aiello[1]

[1] University of Stuttgart, Stuttgart, Germany
{robin.pesl,marco.aiello}@iaas.uni-stuttgart.de
[2] Mercedes-Benz Group AG, Stuttgart, Germany
[3] cellcentric GmbH & Co KG, Kirchheim/Teck-Nabern, Germany

**Abstract.** A classic, central Service-Oriented Computing (SOC) challenge is the service composition problem. It concerns solving a user-defined task by selecting a suitable set of services, possibly found at runtime, determining an invocation order, and handling request and response parameters. The solutions proposed in the past two decades mostly resort to additional formal modeling of the services, leading to extra effort, scalability issues, and overall brittleness. With the rise of Large Language Models (LLMs), it has become feasible to process semi-structured information like state-of-the-practice OpenAPI documentation containing formal parts like endpoints and free-form elements like descriptions. We propose Compositio Prompto to generate service compositions from those semi-structured documentation. Compositio Prompto acts as an encapsulation of the prompt creation and the model invocation such that the user only has to provide the service specifications, the task, and which input and output format they expect, eliminating any manual and laborious annotation or modeling task by relying on already existing documentation. To validate our approach, we implement a fully operational prototype, which operates on a set of OpenAPIs, a plain text task, and an input and output JSON schema as input and returns the generated service composition as executable Python code. We measure the effectiveness of our approach on a parking spot booking case study. Our experiments show that models can solve several tasks, especially those above 70B parameters, but none can fulfill all tasks. Furthermore, compared with manually created sample solutions, the ones generated by LLMs appear to be close approximations.

**Keywords:** Automated service composition · Service discovery · Large language models · Code generation · Automotive services

# Service Composition

# The Case of Service Composition
*Definition*

**Service Composition** is the process of integrating independent loosley coupled services starting from a user request based on the ones available in the execution context. The services communicate over a network and are modular, allowing for flexible and dynamic composition. The *orchestrator* is responsible for coordinating the service composition.

# Service Composition as AI Planning

- **Artificial Intelligence Planning and Scheduling** is a branch of Artificial Intelligence devoted to the study of algorithms and systems to empower intelligent agents with the ability to pursue their goals.

- *Goal:* a description of the state of the world to realise
  ### *user request*

- *Plan:* an algorithm that describes how to reach a goal state
  ### *a composition to orchestrate*

- *Environment:* a system the state of which can be sensed and changed by the planning actor
  ### *APIs, service states, domain knowledge*

# Large Reasoning Models

- A neural network-based model optimized for multi-step logical and symbolic reasoning

- Trained on heterogeneous datasets: natural language, formal logic, math, code, and multimodal inputs

- Excels at structured problem-solving via in-context learning, chain-of-thought prompting, and tool augmentation

- Designed to perform algorithmic reasoning, planning, and hypothetical simulation

- May incorporate external memory, RAG, or tool use (e.g., calculators, search APIs, WolframAlpha)

- Good for: automated theorem proving, scientific discovery, decision support, etc.

# The Illusion of Thinking:
# Understanding the Strengths and Limitations of Reasoning Models via the Lens of Problem Complexity

Parshin Shojaee[*][†]    Iman Mirzadeh[*]    Keivan Alizadeh
Maxwell Horton    Samy Bengio    Mehrdad Farajtabar

Apple

**Abstract**

Recent generations of frontier language models have introduced Large Reasoning Models (LRMs) that generate detailed thinking processes before providing answers. While these models demonstrate improved performance on reasoning benchmarks, their fundamental capabilities, scaling properties, and limitations remain insufficiently understood. Current evaluations primarily focus on established mathematical and coding benchmarks, emphasizing final answer accuracy. However, this evaluation paradigm often suffers from data contamination and does not provide insights into the reasoning traces' structure and quality. In this work, we systematically investigate these gaps with the help of controllable puzzle environments that allow precise manipulation of compositional complexity while maintaining consistent logical structures. This setup enables the analysis of not only final answers but also the internal reasoning traces, offering insights into how LRMs "think". Through extensive experimentation across diverse puzzles, we show that frontier LRMs face a complete accuracy collapse beyond certain complexities. Moreover, they exhibit a counter-intuitive scaling limit: their reasoning effort increases with problem complexity up to a point, then declines despite having an adequate token budget. By comparing LRMs with their standard LLM counterparts under equivalent inference compute, we identify three performance regimes: (1) low-complexity tasks where standard models surprisingly outperform LRMs, (2) medium-complexity tasks where additional thinking in LRMs demonstrates advantage, and (3) high-complexity tasks where both models experience complete collapse. We found that LRMs have limitations in exact computation: they fail to use explicit algorithms and reason inconsistently across puzzles. We also investigate the reasoning traces in more depth, studying the patterns of explored solutions and analyzing the models' computational behavior, shedding light on their strengths, limitations, and ultimately raising crucial questions about their true reasoning capabilities.

# STOP ANTHROPOMORPHIZING INTERMEDIATE TOKENS AS REASONING/THINKING TRACES!

**Subbarao Kambhampati**    **Kaya Stechly**    **Karthik Valmeekam**    **Lucas Saldyt**    **Siddhant Bhambri**

**Vardhan Palod**    **Atharva Gundawar**    **Soumya Rani Samineni**    **Durgesh Kalwar**    **Upasana Biswas**

**School of Computing & AI**
**Arizona State University**

## ABSTRACT

Intermediate token generation (ITG), where a model produces output before the solution, has been proposed as a method to improve the performance of language models on reasoning tasks. These intermediate tokens have been called "reasoning traces" or even "thoughts" – implicitly anthropomorphizing the model, implying these tokens resemble steps a human might take when solving a challenging problem. In this paper, we present evidence that this anthropomorphization isn't a harmless metaphor, and instead is quite dangerous – it confuses the nature of these models and how to use them effectively, and leads to questionable research.

v2 [cs.AI] 27 May 2025

# Large Action Models

- A parameter-rich neural policy model trained to map from high-dimensional observations and goals to action distributions

- Leverages transformer-based architectures for sequence modeling of action trajectories

- Operates over state-action-return triples (or variations) for temporal credit assignment and long-horizon planning

- Trained via offline reinforcement learning, behavior cloning, or trajectory-level supervision from expert demonstrations or synthetic data

- Can ingest multimodal inputs and output low-level control signals or symbolic action commands

- Supports zero-shot generalization across tasks via goal-conditioning, prompting, or language grounding

- Frequently deployed in embodied agents, robotic manipulation, navigation, game environments, and tool use contexts

# Language Models (LMs)

## Large Language Models (LLMs)

- Large-scale knowledge of language patterns
- Understands natural language
- Excels at generating coherent text
- Efficient for translation Q&A, etc.

- No deliberate, iterative reasoning
- No direct interaction with environment
- Tends to hallucinate

GPT-3.5, GPT-4, Claude, LLaMa, BERT, Qwen, Grok, Gemini 1

## Large Reasoning Models (LRMs)

- Understands natural language
- Deep, explicit reasoning
- Strong at planning and problem-solving
- Often uses explicit structure of reasoning

- Typically does not act on the environment
- May require extensive computational resources
- Can be slower

OpenAI o1/o3, DeepSeek R1, Gemini 2.0, QwQ

## Large Action Models (LAMs)

- Understands multimodal inputs
- Directly executes actions in environments
- Integrates sensing with action outputs
- Can handle tasks in real time

- Often weak at (high-level) reasoning
- Requires specialised data (e.g., action logs)
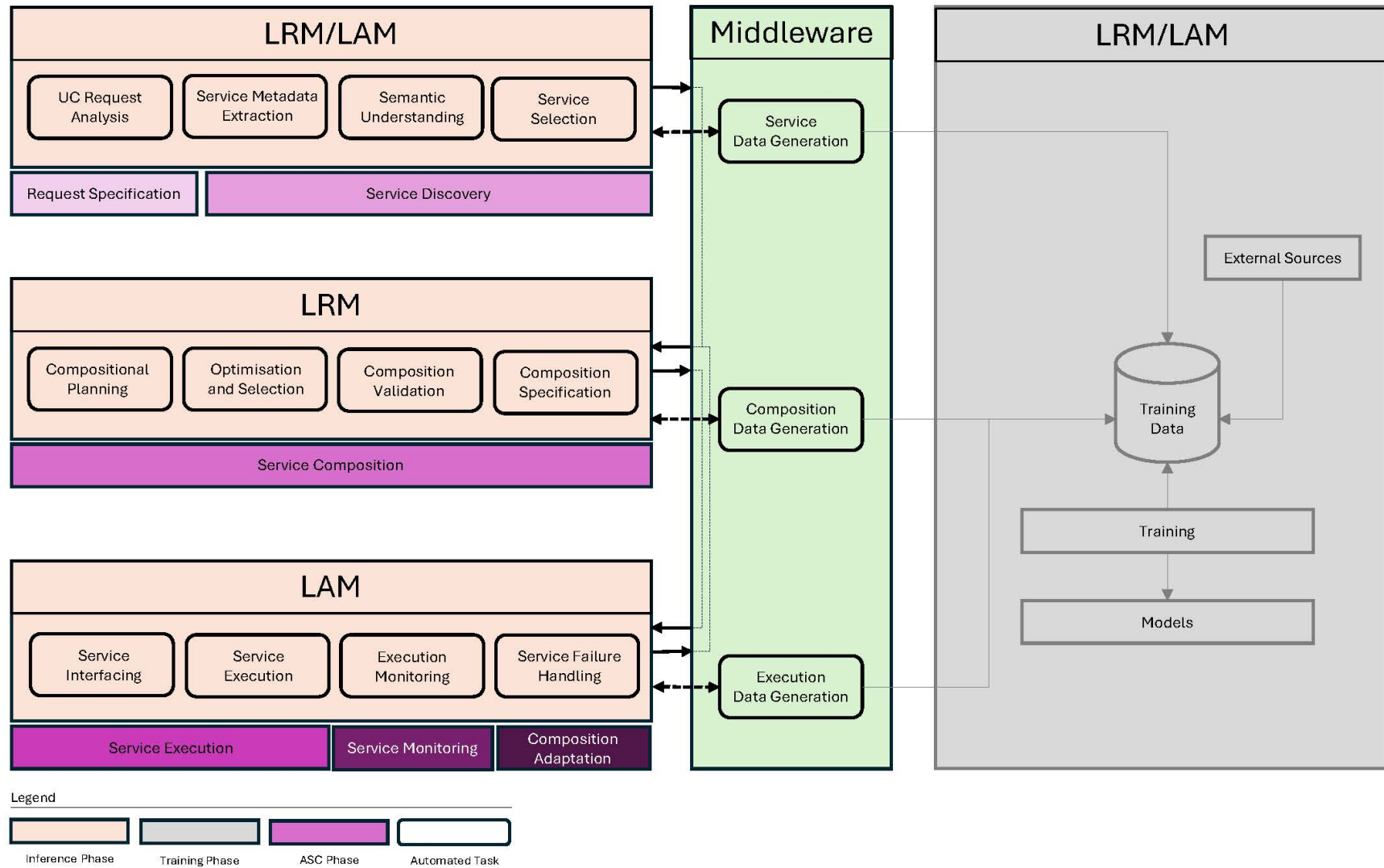- Can be difficult to train for safety and reliability

Google RT-1/RT-2, DeepMind Gato, Rabbit R1, CogAgent, ScreenAI, xLAM

Legend

| Strengths | Limitations | Examples |

**LRM/LAM**

| UC Request Analysis | Service Metadata Extraction | Semantic Understanding | Service Selection |

Request Specification | Service Discovery

**Middleware**

Service Data Generation

Composition Data Generation

Execution Data Generation

**LRM/LAM**

External Sources

Training Data

Training

Models

**LRM**

| Compositional Planning | Optimisation and Selection | Composition Validation | Composition Specification |

Service Composition

**LAM**

| Service Interfacing | Service Execution | Execution Monitoring | Service Failure Handling |

Service Execution | Service Monitoring | Composition Adaptation

Legend

Inference Phase | Training Phase | ASC Phase | Automated Task
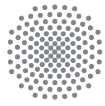
©Marco Aiello, 2025

# Concluding remarks

# Reflection on Initial Steps

- LLM, LRM, LAM can cover various aspects of Service Composition

- Promising technologies with some known and yet unknown limitations

- See you at SummerSOC 2026 for more

**Universität Stuttgart**
Institute of Architecture of Application Systems

# Thank you!   Vielen Dank!   Grazie!   Merci!   Bedankt!

**Marco Aiello**

多謝!

E-Mail   marco.aiello@iaas.uni-stuttgart.de

Telefon  +49 (0) 711 685-88471

www.iaas.uni-stuttgart.de/en/department-service-computing/  or aiellom.it

Universität Stuttgart

Service Computing Dep  — IAAS

Universitätsstr. 38, 70569 Stuttgart

Questions?