

Real-world Al-based Systems

Christian Cabrera Jojoa Research Associate chc79@cam.ac.uk

Department of Computer Science and Technology

Agenda

PART I: The great AI question.

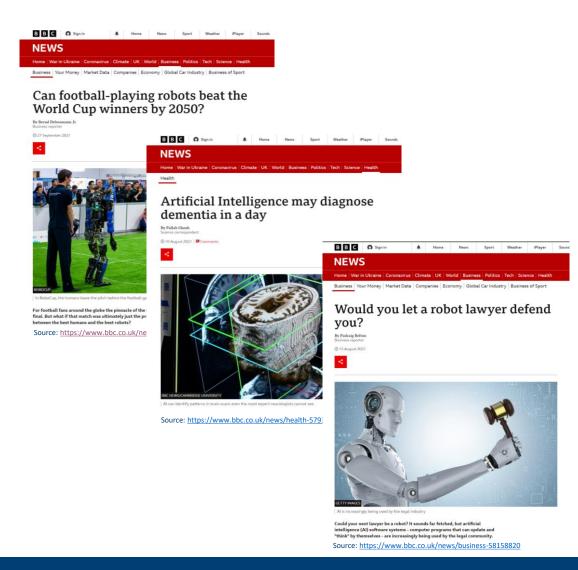
PART II: Why the question?

PART III: Possible answers...

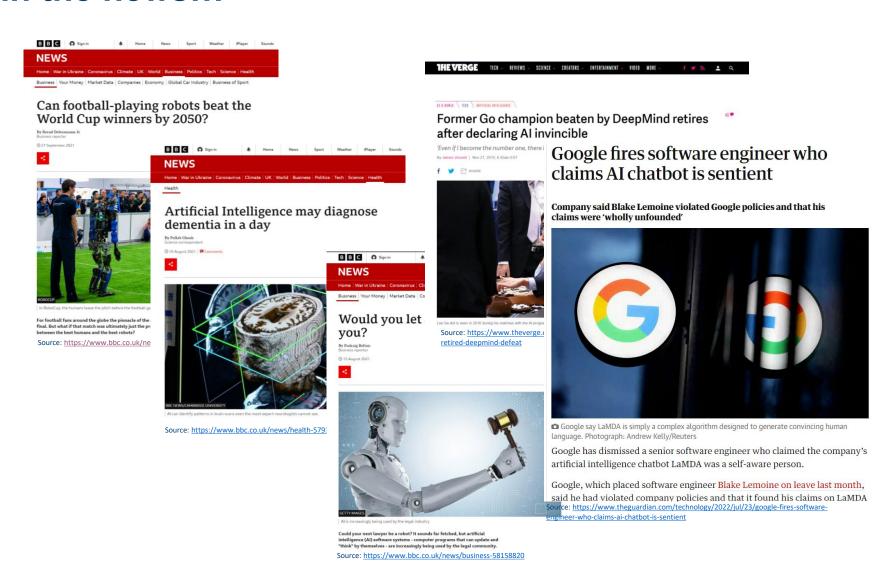


PART I: The great AI question.

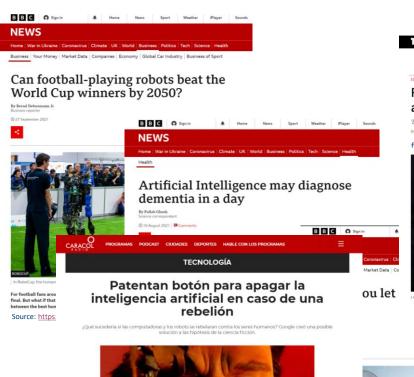


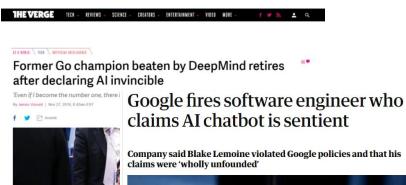












Source: https://www.theverge.org retired-deepmind-defeat



Source: https://caracol.com.co/radio/2016/06/15/tecnologia/1466026084 693727.htmlg/poi

Could your next lawyer be a robot? It sounds far fetched, but artificial intelligence (AI) software systems - computer programs that can update and "think" by themselves - are increasingly being used by the legal community. Source: https://www.bbc.co.uk/news/business-58158820

21.3% Growth Forecasted For Artificial Intelligence (AI) Edge Computing Market, Surpassing \$53.18 Billion By 2029



2025 Market Reports Update: Forecasts Through 2034, Emerging Trends, Key Players, and Leading Regions - Stay Ahead of the Competition

> Source: https://www.whatech.com/og/markets-research/it/957645-21-3-growth-forecasted-for-artificial-intelligence-ai-edge-computing market-surpassing-53-18-billion-by-2029.html

6AM HOY POR HOY

La inteligencia artificial lo cambiará todo

Escuche el análisis de Juan Carlos Echeverry sobre este tema que se está tomando el mundo

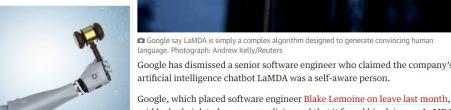


LA INTELIGENCIA ARTIFICIAL LO CAMBIARÁ TODO



https://caracol.com.co/programa/2019/03/2 1/6am_hoy_por_hoy/1552309220_833279.l

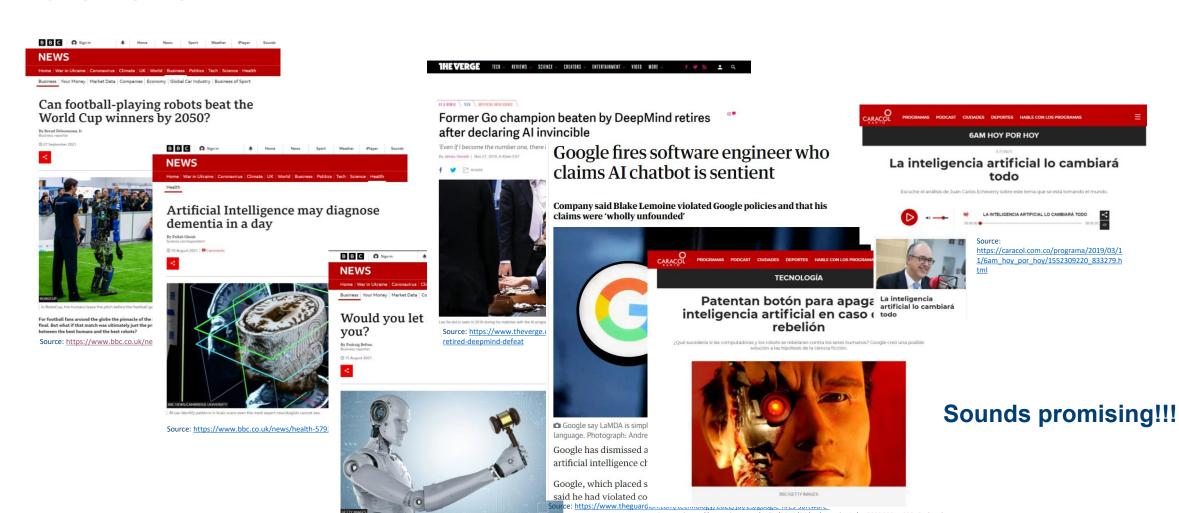






said he had violated company policies and that it found his claims on LaMDA Source: https://www.theguardian.com/technology/2022/jul/23/google-fires-softwarengineer-who-claims-ai-chatbot-is-sentient





neer-who-claims-ai-chatbot-isounce:rhttps://caracol.com.co/radio/2016/06/15/tecnologia/1466026084 693727.html

Could your next lawyer be a robot? It sounds far fetched, but artificial intelligence (AI) software systems - computer programs that can update and "think" by themselves - are increasingly being used by the legal community.

Source: https://www.bbc.co.uk/news/business-58158820

UNIVERSITY OF CAMBRIDGE

But...



TRANSPORTATION | TESTA | FEATURED STORIES | Source: https://www.theverge.com/2021/7/5/22563751/tesla-elon-musk-full-self-driving-admission-autopilot-crash

Elon Musk just now realizing that self-driving cars are a 'hard problem'



Overpromising and underdelivering

By Andrew J. Hawkins | @andyjayhawk | Jul 5, 2021, 10:53am EDT

6AM HOY POR HOY La inteligencia artificial lo cambiará todo "Twilio is ranked #1 Escuche el análisis de Juan Carlos Echeverry sobre este tema que se está tomando el mundo. for market share in Customer Data Platfo N LA INTELIGENCIA ARTIFICIAL LO CAMBIARÁ TODO IDC, "Worldwide Customer Dat Platform Market Shares, 2020 Source: https://caracol.com.co/programa/2019/03/1 1/6am_hov_por_hoy/1552309220_833279.h t TECNOLOGÍA otón para apaga La inteligencia artificial lo cambiará rtificial en caso (todo rebelión verge deáls Subscribe to get the best Ve approved tech deals of the Email (required) It is far from reality... By signing up, you agree to our Privac and European users agree to the data SUBSCRIBE BBC/GETTY IMAGES

dio/2016/06/15/tecnologia/1466026084 693727.html

Tesla CEO Elon Musk is finally admitting that he underestimated how difficult it is to develop a safe and reliable self-driving car. To which the entire engineering community rose up as one to say, "No duh."

Could your next lawyer be a robot? It sounds far fetched, but artificial intelligence (Al) software systems - computer programs that can update and "think" by themselves - are increasingly being used by the legal community.

Source: https://www.bbc.co.uk/news/business-58158820



The Great Al Fallacy

Al in the news...



"The fallacy is associated with an implicit promise that is embedded in many statements about Artificial Intelligence. Artificial Intelligence, as it currently exists, is merely a form of automated decision making. The implicit promise of Artificial Intelligence is that it will be the first wave of automation where the machine adapts to the human, rather than the human adapting to the machine."[1]

Neil Lawrence
DeepMind Professor
of Machine Learning
at the University
of Cambridge

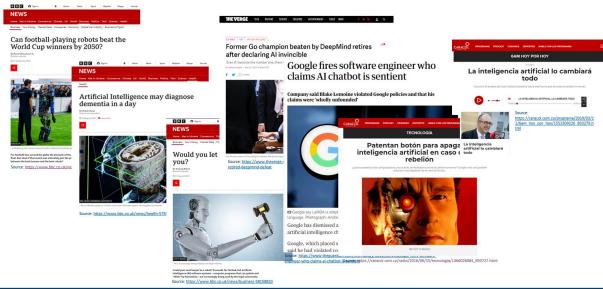


[1] Lawrence N. "The Great AI Fallacy". Webinar for The Cambridge Network, Apr 2020, Available online: http://inverseprobability.com/talks/notes/the-great-ai-fallacy.html



The Great AI Question

Al in the news...

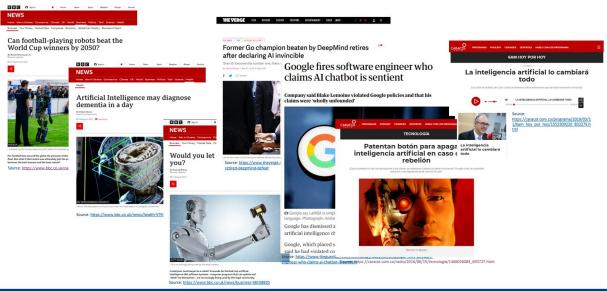






The Great AI Question

Al in the news...



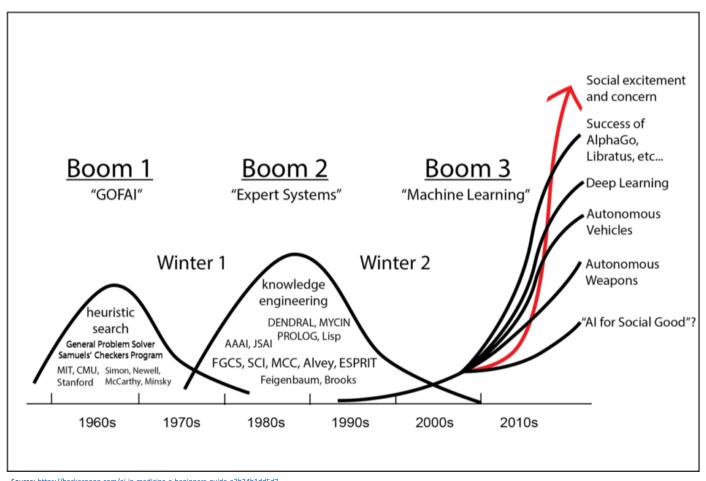
UNIVERSITY OF CAMBRIDGE

- The News focus on impressive headlines.
- But they lack scientific rigor:
 - Problems are not well explained.
 - Assumptions are not considered.
 - Experimental setups are ignored.
 - Limitations are not discussed.
- Overwhelmed optimism!



Overwhelmed optimism

- New technologies usually create this optimism:
 - There were predictions about Humans conquering other worlds by first decades of 2000s.
- Then, we have had disappointment periods:
 - Spatial programs are not as attractive and popular as before.
- Then, confidence is back again to evaluate actual results:
 - Global Position System (GPS) is a tool we all use everyday!



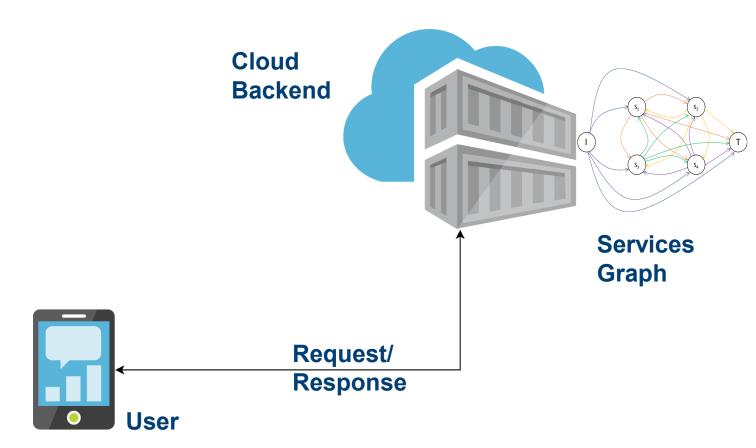
Source: https://hackernoon.com/ai-in-medicine-a-beginners-guide-a3b34b1dd5d7



PART II: Why the question?

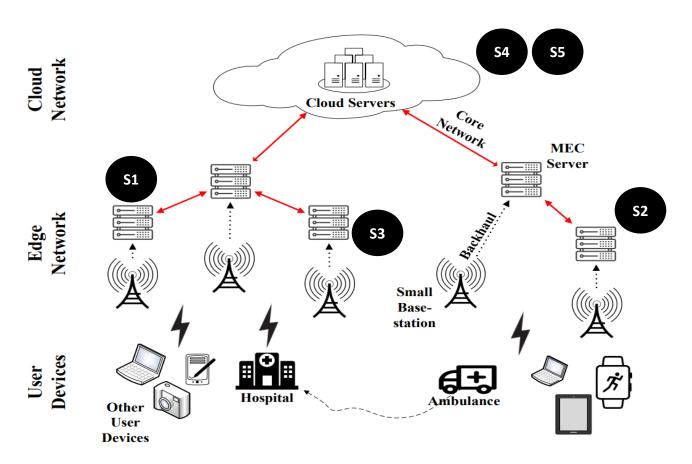


- Current applications are cloud based and they follow a Microservices architecture.
- There is a physical distance between users and the cloud, which generates additional latency.
- Some applications do not tolerate such latency.
- Edge computing is a new paradigm that enables low-latency applications[1].



[1] Malazi, H. T., Chaudhry, S. R., Kazmi, A., Palade, A., Cabrera, C., White, G., & Clarke, S. (2022). Dynamic Service Placement in Multi-access Edge Computing: a Systematic Literature Review. *IEEE Access*.





- **Edge servers** are closed to end users, and they can **process data locally.**
- Services are executed on edge servers, but they have **limited resources**.
- The problem is to find the optimal combination of services and edge servers that minimizes latency subject to available resources.
- It is known as Service Placement Problem[1].

[1] Malazi, H. T., Chaudhry, S. R., Kazmi, A., Palade, A., Cabrera, C., White, G., & Clarke, S. (2022). Dynamic Service Placement in Multi-access Edge Computing: a Systematic Literature Review. *IEEE Access*.



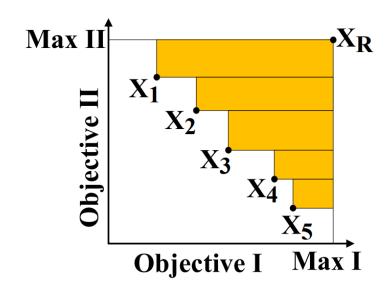
$$\min_{j} wt_j + lat(D_j) (1)$$

$$\min_{i} \sqrt{w_p \sigma(RR^{CPU})^2 + w_m \sigma(RR^{RAM})^2}$$
 (2)

subject to:
$$\sum_{\substack{k=1\\n}}^{n} rr_{ik}^{CPU} <= r_i^{CPU} \quad (3)$$

$$\sum_{k=1}^{n} rr_{ik}^{RAM} <= r_i^{RAM} \qquad (4)$$

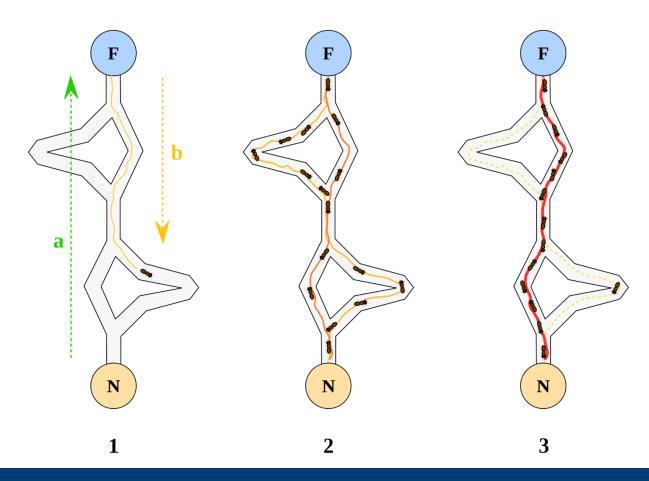
Pareto-optimal[1]



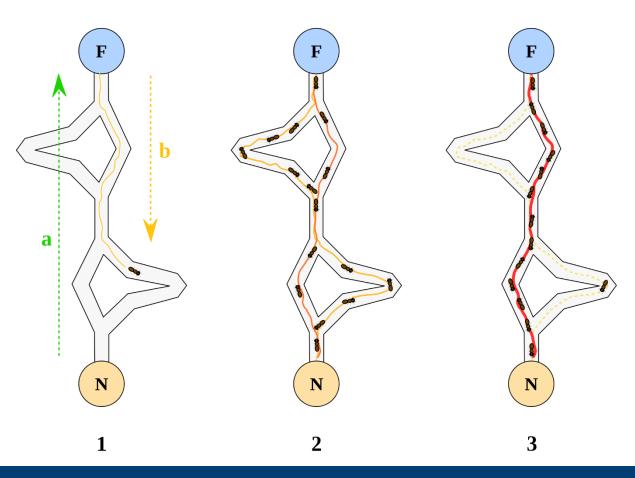
[1] Palade, A., Cabrera, C., White, G., & Clarke, S. (2018, November). Stigmergic service composition and adaptation in mobile environments. In *International Conference on Service-Oriented Computing* (pp. 618-633). Springer, Cham.



- Ant-Colony Optimisation



Ant-Colony Optimisation

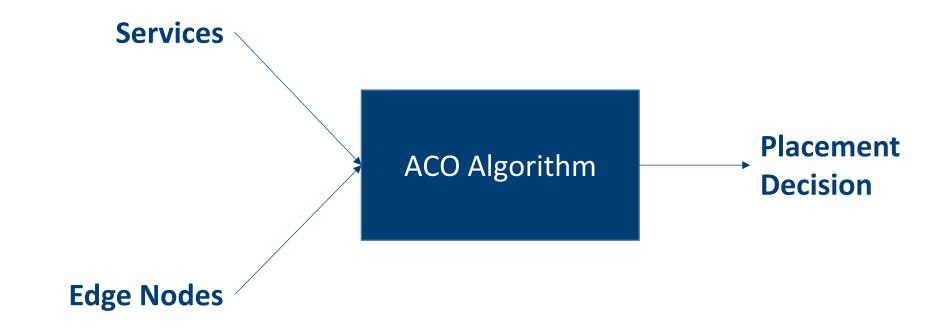


Forward movement:

$$p(t,i) = \frac{\left[\tau_{t,i}\right]^{\alpha} * \left[\eta_{i,i}\right]^{\beta}}{\sum_{i=1}^{n} \left[\tau_{t,i}\right]^{\alpha} * \left[\eta_{j,i}\right]^{\beta}}$$

Backward movement:

$$\tau_{t,i} = \tau_{t,i} + \boldsymbol{\rho}$$



 $App_n = < S, RQ, I, p >$

$$N_1$$
 $N_1 = \langle R, L, I \rangle$

$$N_2$$
 $N_2 = \langle R, L, I \rangle$

$$N_3$$
 $N_3 = \langle R, L, I \rangle$

App₁=<S, RQ, I, p> S=< s_1 , s_2 , s_3 > RQ=< R_{f1} , R_{f2} , Rf_{f3} >









$$N_2$$

$$N_2$$



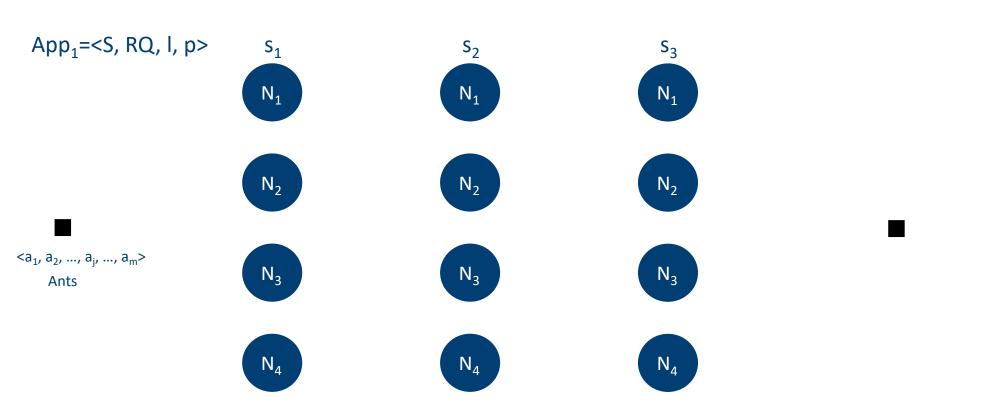
$$N_3$$

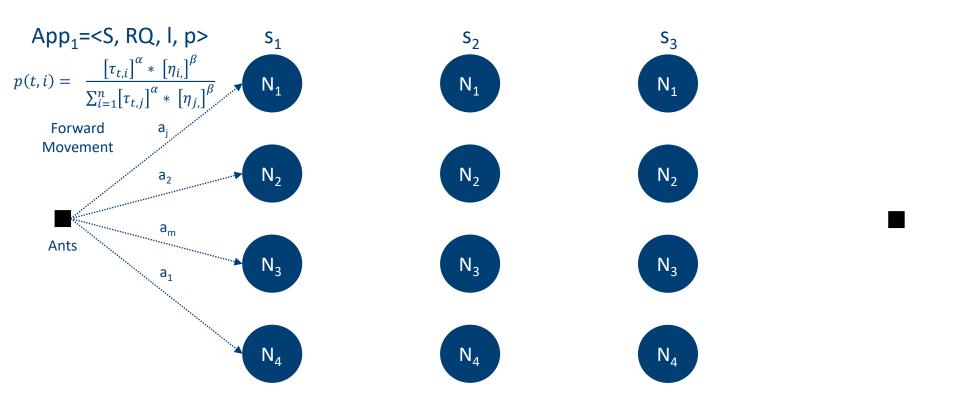
$$N_3$$

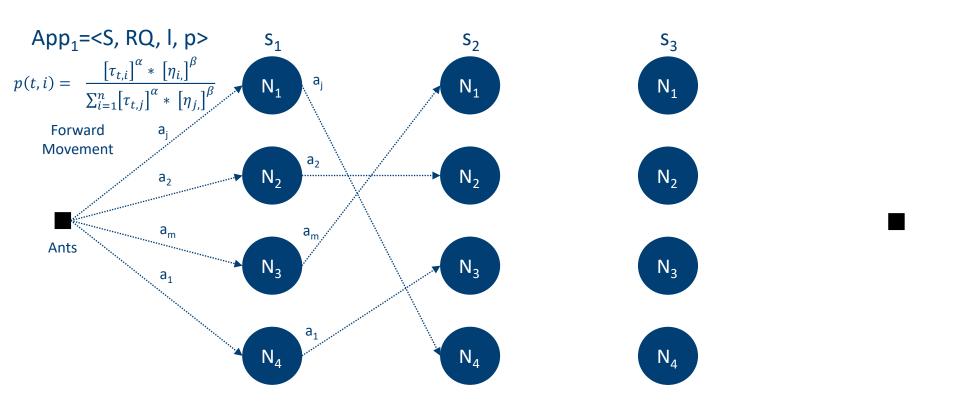
$$N_4$$

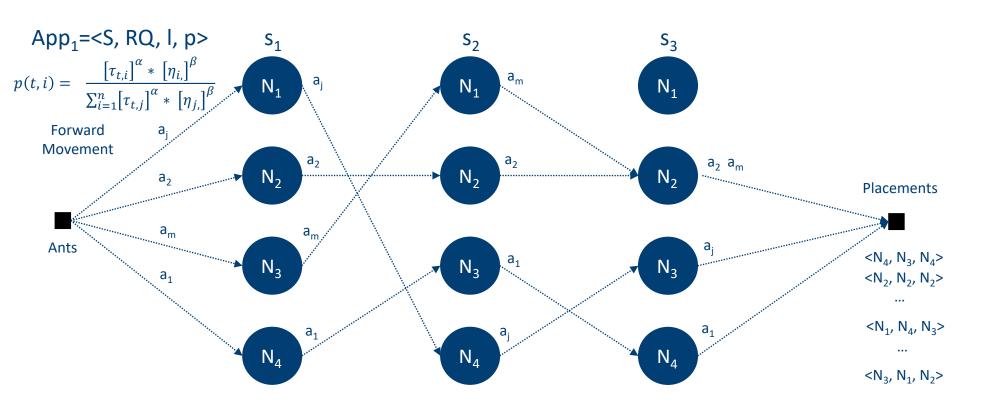
$$N_4$$

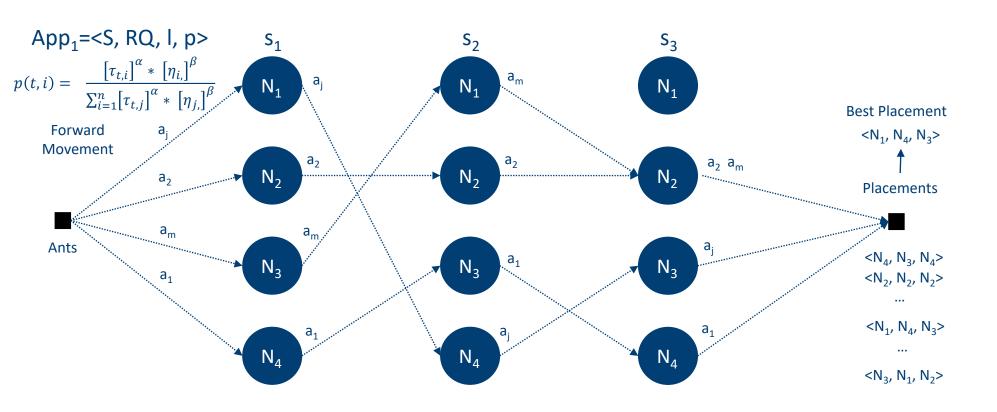
$$N_4$$

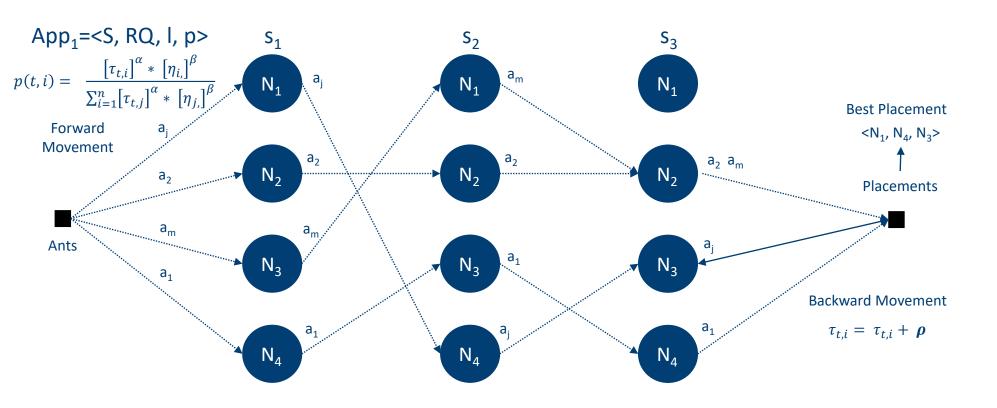


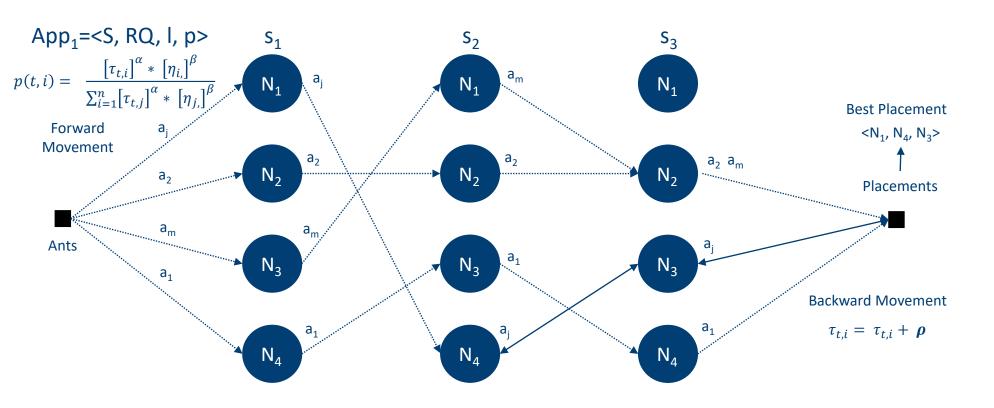


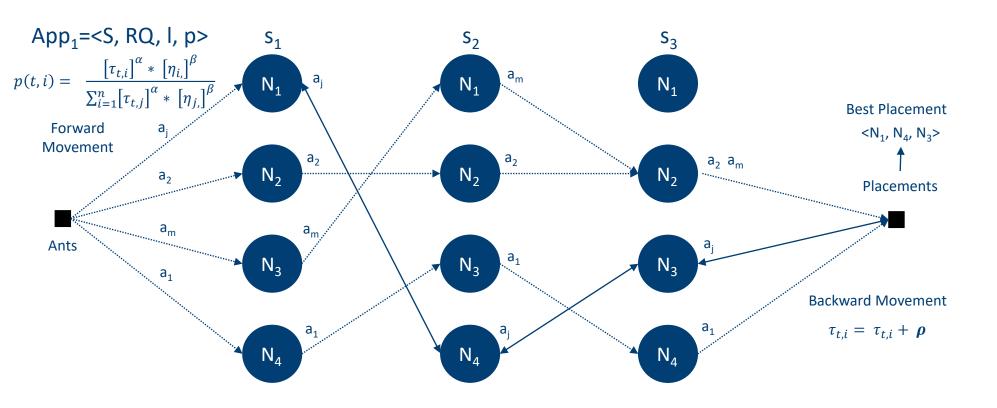


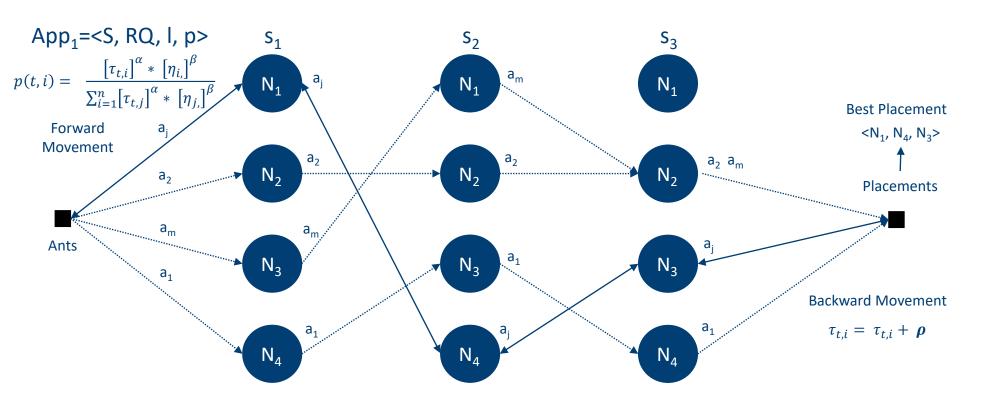


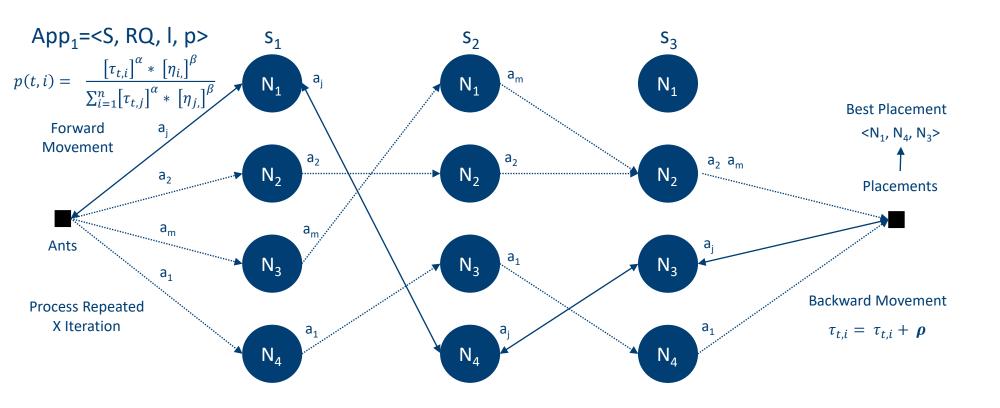


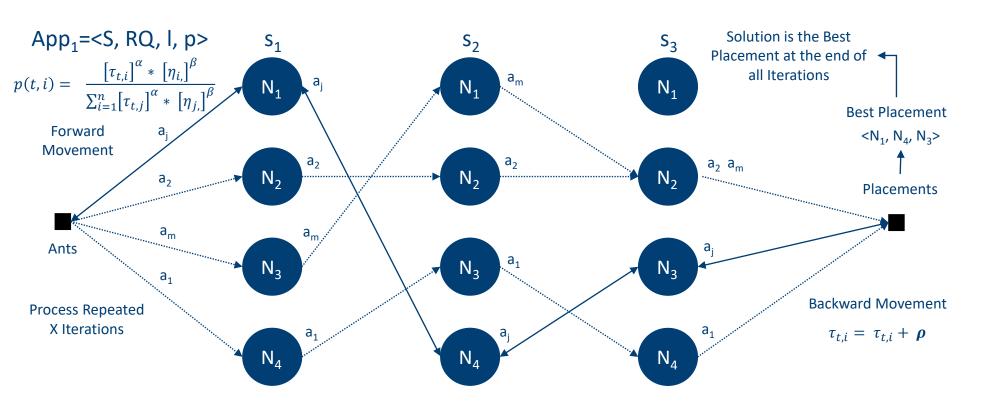


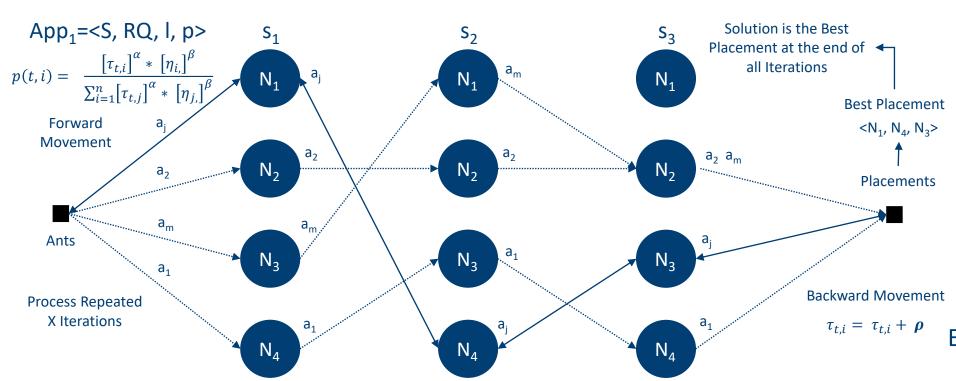










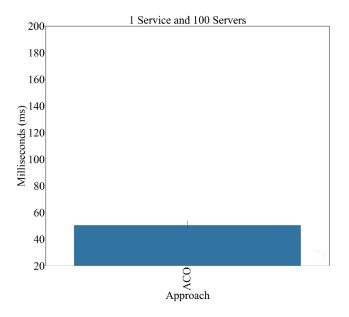


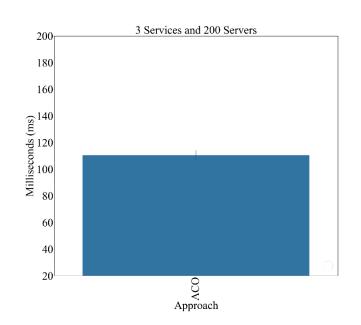
Execution time is affected by:

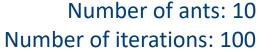
- Number of services.
- Number of edge servers.
- Number of ants.
- Number of iterations

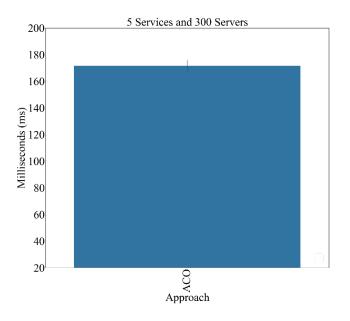


Smart-city simulation[1]



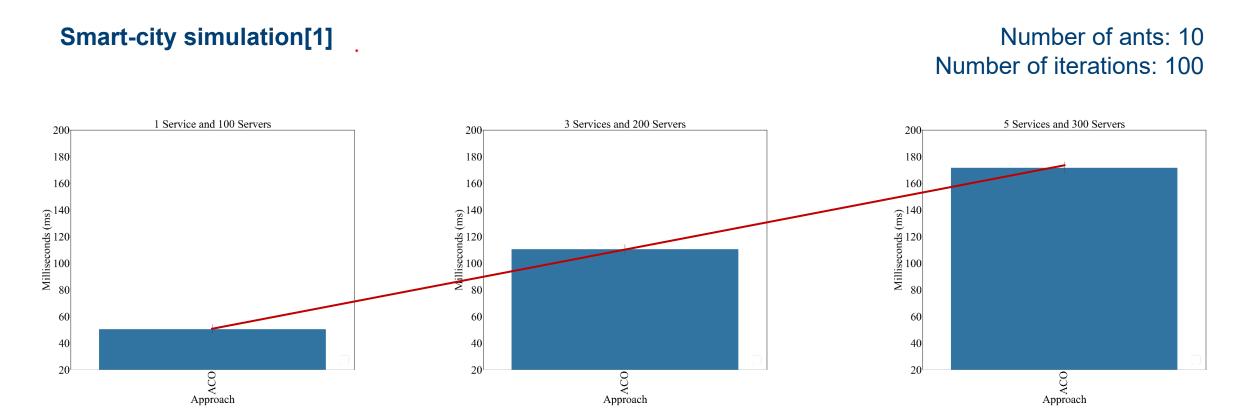






[1] Richerzhagen, B., Stingl, D., Ruckert, J., & Steinmetz, R. (2015, August). Simonstrator: Simulation and prototyping platform for distributed mobile applications. In *The 8th EAI International Conference on Simulation Tools and Techniques (ACM SIMUTOOLS 2015)* (pp. 99-108).

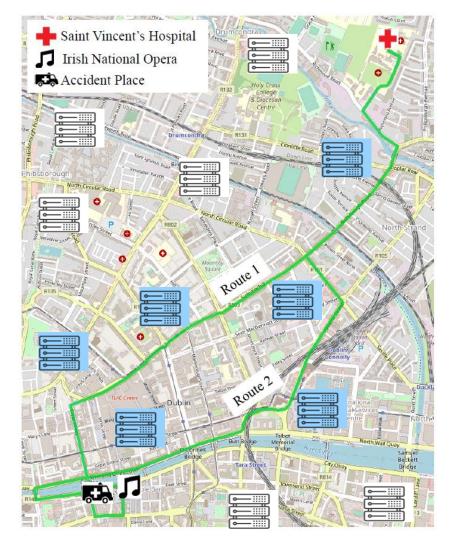




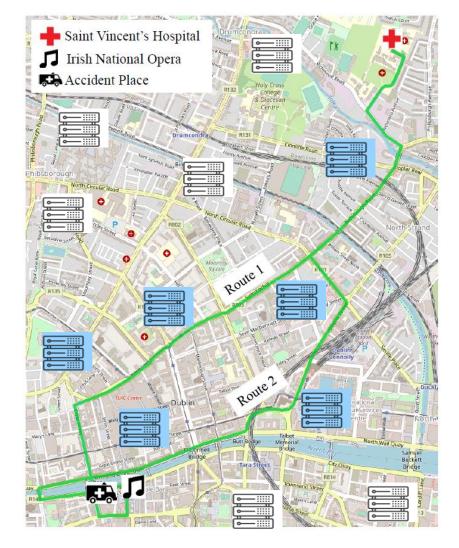
This execution time does not suit low-latency requirements, but that is how ACO is designed!!!

[1] Richerzhagen, B., Stingl, D., Ruckert, J., & Steinmetz, R. (2015, August). Simonstrator: Simulation and prototyping platform for distributed mobile applications. In *The 8th EAI International Conference on Simulation Tools and Techniques (ACM SIMUTOOLS 2015)* (pp. 99-108).



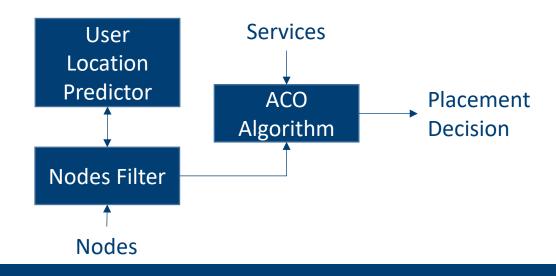


- Variables we cannot reduce:
 - Number of services.
 - Number of iterations.
 - Number of ants.
- We can reduce the number of servers:
 - We can pre-select them if we can predict user locations.



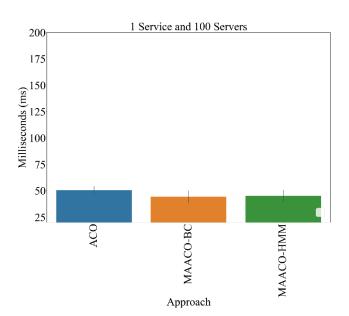


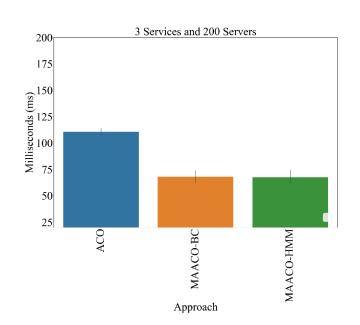
Re-design of the solution:



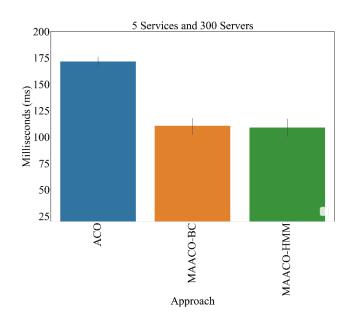


Smart-city simulation[1]









[1] Richerzhagen, B., Stingl, D., Ruckert, J., & Steinmetz, R. (2015, August). Simonstrator: Simulation and prototyping platform for distributed mobile applications. In *The 8th EAI International Conference on Simulation Tools and Techniques (ACM SIMUTOOLS 2015)* (pp. 99-108).



Bayesian Classifier [1]

- Transition matrix depends on the number of links, which are the number of streets in a city.
- A lot of data (i.e., trips) are needed to train the model.
- Training time is now an issue!
- We assumed a limited number of streets in our work.

Hidden Markov Model [2]

- Frequency matrix depends on the number of links, which are the number of streets in a city.
- A lot of data (i.e., trips) are needed to train the model.
- Training time is now an issue!
- We assumed a limited number of streets in our work.

^[2] Y. Lassoued, J. Monteil, Y. Gu, G. Russo, R. Shorten, and M. Mevissen, "A hidden markov model for route and destination prediction," n 2017 IEEE 20th International Conference on Intelligent Transportation Systems (ITSC). IEEE, 2017.



^[1] J. P. Epperlein, J. Monteil, M. Liu, Y. Gu, S. Zhuk, and R. Shorten, "Bayesian classifier for route prediction with markov chains," in 2018 21st International Conference on Intelligent Transportation Systems (ITSC). IEEE, 2018, pp. 677–682.

Bayesian Classifier [1]

- Transition matrix depends on the number of links, which are the number of streets in a city.
- A lot of data (i.e., trips) are needed to train the model.
- Training time is now an issue!
- We assumed a limited number of streets in our work.

Hidden Markov Model [2]

- Frequency matrix depends on the number of links, which are the number of streets in a city.
- A lot of data (i.e., trips) are needed to train the model.
- Training time is now an issue!
- We assumed a limited number of streets in our work.

AGAIN, NEW DESIGN DECISIONS ARE NEEDED TO DEPLOY THESE ML ALGORITHMS IN THE REAL-WORLD!

[1] J. P. Epperlein, J. Monteil, M. Liu, Y. Gu, S. Zhuk, and R. Shorten, "Bayesian classifier for route prediction with markov chains," in 2018 21st International Conference on Intelligent Transportation Systems (ITSC). IEEE, 2018, pp. 677–682.

[2] Y. Lassoued, J. Monteil, Y. Gu, G. Russo, R. Shorten, and M. Mevissen, "A hidden markov model for route and destination prediction," n 2017 IEEE 20th International Conference on Intelligent Transportation Systems (ITSC). IEEE, 2017.



Al Deployment in the Real-world

- Real-world does not support assumptions.
- Real world environments are usually large, heterogeneous, complex and dynamic.

How to support engineers and developers in the process of deploying Al-based systems in real-world environments?

How to enable the reasoning about Al-based systems in real world environments?

How to automate the repetitive design process?



PART III: Possible answers...



The AutoAl Programme scales our ability to deploy safe and reliable Al solutions, driving innovation in machine learning-enabled techniques for deploying, maintaining, and understanding Al systems. By investigating how to decompose Al systems into their component parts, how to manage data in system development, and how to monitor performance in deployment, AutoAl will develop a new Al design and engineering paradigm.





faculty









faculty







17 4







facu

PROBLEM FIRST!!!















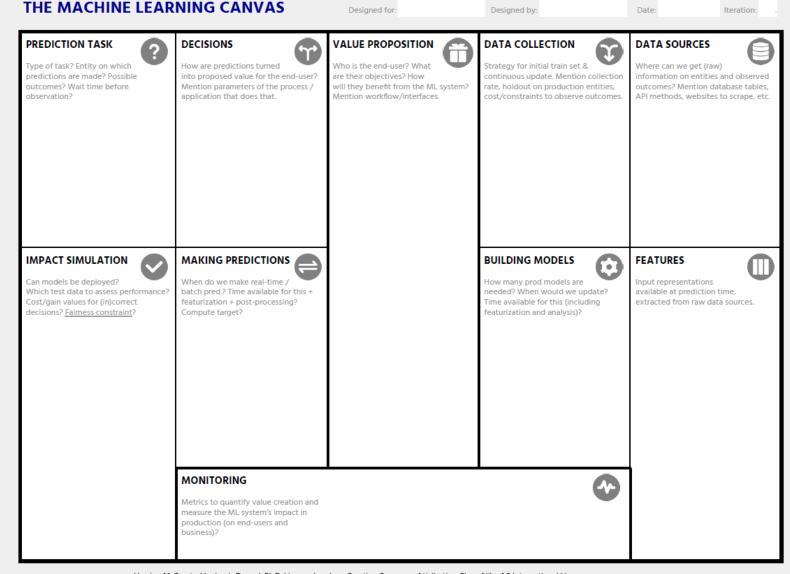








- What is the problem?
- Functional vs Non-functional requirements?
- Why do we need Al?
- How would AI impact our requirements?





Version 1.1. Created by Louis Dorard, Ph.D. Licensed under a <u>Creative Commons Attribution-ShareAlike 4.0 International License</u>. Please keep this mention and the link to <u>ownmLco</u> when sharing.

OWNML.CO



compute $model + data \rightarrow prediction$

compute $model + data \rightarrow prediction$

- There has been **a lot** of focus on and research **efforts** in learning **models** (i.e., new algorithms) and **compute** (i.e., more powerful machines).
- But, the data has not gotten enough attention!

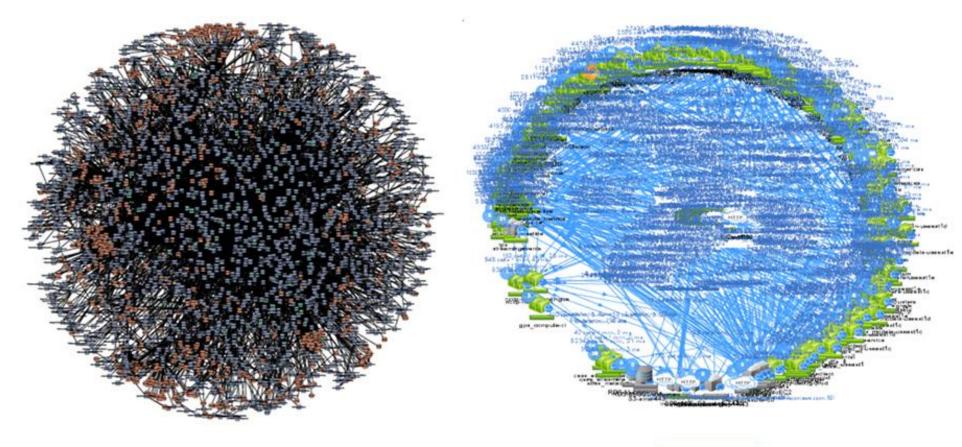
compute $model + data \rightarrow prediction$

- There has been **a lot** of focus on and research **efforts** in learning **models** (i.e., new algorithms) and **compute** (i.e., more powerful machines).
- But, the data has not gotten enough attention!

The Data Dichotomy:

"While data-driven systems are about exposing data, service oriented architectures are about hiding data."





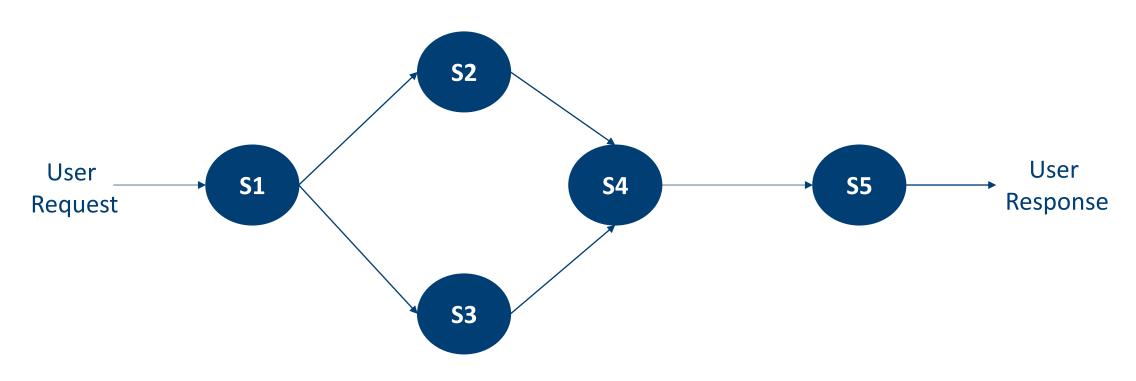




Source: https://www.divante.com/blog/10-companies-that-implemented-the-microservice-architecture-and-paved-the-way-for-others

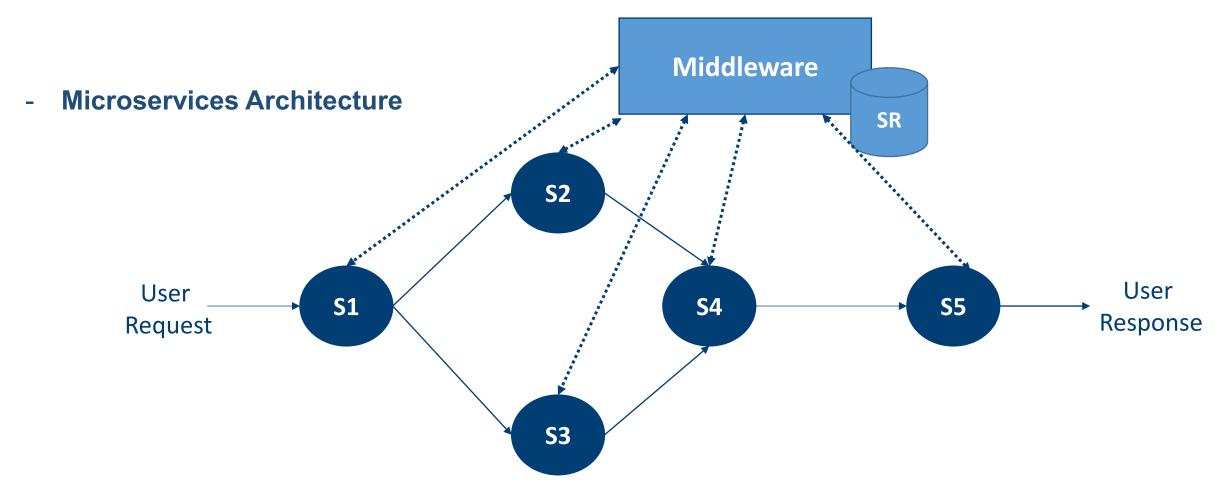


- Microservices Architecture



Operations are first!





Operations are first!



Data Oriented Architecture

S2

User Request/Response **S4**

Publish/Subscribe

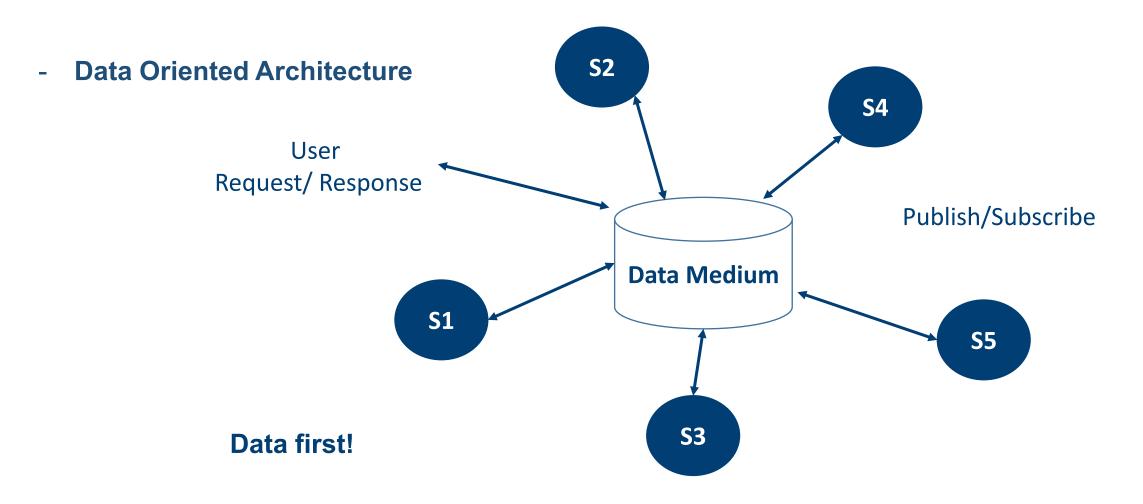
S1

S5

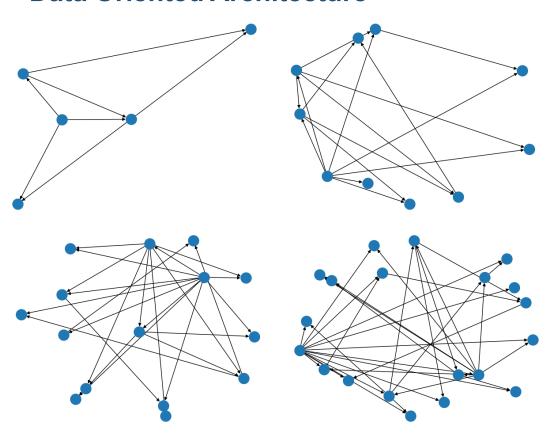
Data first!

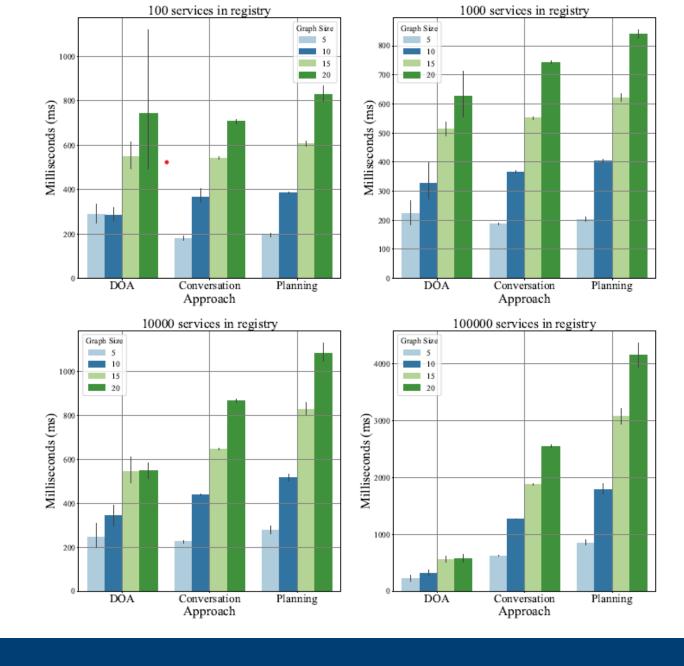
S3





Data Oriented Architecture







Take aways...

- It is very likely most of the AI popular predictions will not become real.
- Al has potential but we must be careful with the overwhelmed optimism.
- We must approach AI progress in a scientific way.
- Problems must be always first!
- The deployment of Al-based systems in real-world environments is hard. (**Design Decisions**)
- New paradigms are needed to realize the potential of AI. (Data Oriented Architectures)



Thank you!

chc79@cam.ac.uk

