

SIMULATING LLM ECOSYSTEMS

FROM REFERENCE ARCHITECTURES OF LLM ECOSYSTEMS TO
SIMULATING PERFORMANCE, SUSTAINABILITY, EFFICIENCY


@Large Research
Massivizing Computer Systems

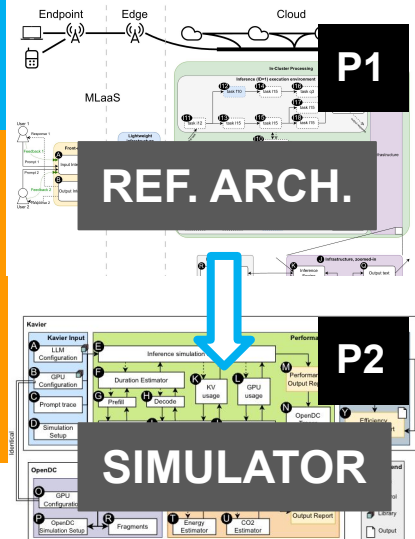
These slides:
bit.ly/radu-bsc

Towards digital twinning
LLM Ecosystems
= **ref. arch.** + **simulator** +
monitoring, datagen + ODA
+goal-oriented infr. steering



Radu Nicolae
@VU Amsterdam
@Large Research

 @rnicolae,
mail@radu-nicolae.com



LLM(s)

(ecosystems under inference)

LLMs are widely adopted everywhere


LLMs are being widely adopted, worldwide, amongst academia, government, and industry.

(Microsoft)scaling at all costs



Source: The guardian ([online](#))

SUSTAINABILITY

Inference GPT: 1 GWh daily
Equivalent to ~40,000 dutch 

FINANCIAL

Inference GPT: \$700k daily
(operational costs)

PERFORMANCE

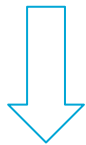
Modern day moore's law

Deploying, monitoring, and scaling LLMs

LLMs

VERY LARGE

VERY COSTLY



WE NEED TO SIMULATE FIRST

Simulators enable performing large-scale experiments in a time and cost efficient way.

Challenge:

No simulator for predicting LLM performance & sustainability

P2

Approach:

Build sim. adhering to ref. arch.

Challenge:

No (good) reference architecture of LLM ecosys under inference

P1

Tackling:

Design+validate ref arch.

A REFERENCE ARCHITECTURE FOR LLM ECOSYSTEMS UNDER INFERENCE

DESIGN [1/4]

Design requirements^[1]

Validity

Usefulness

Design principles

DP1 Heterogeneity

DP2 End-to-end workflow

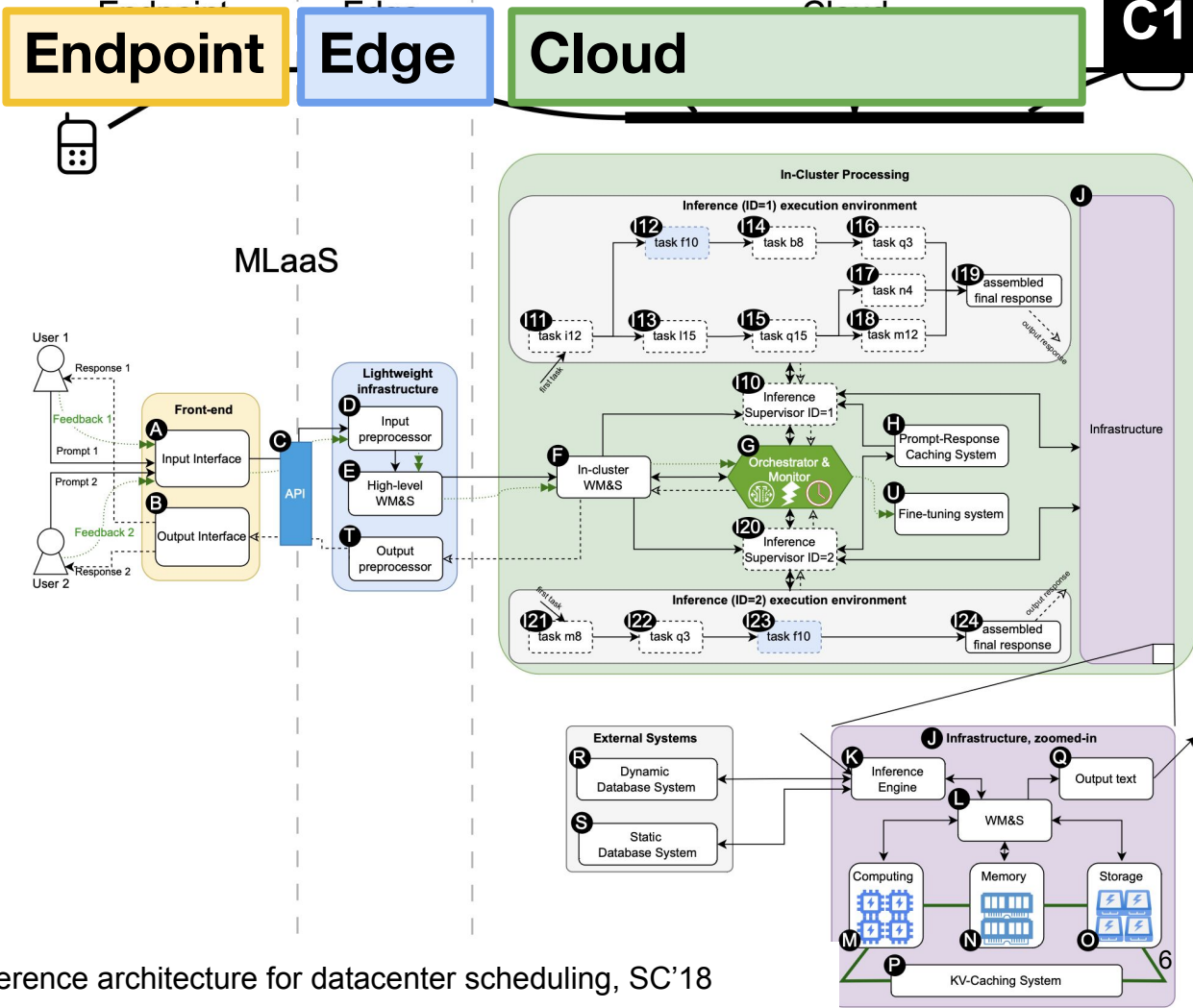
DP3 Branching prompts

DP4 Multi-user ecosys

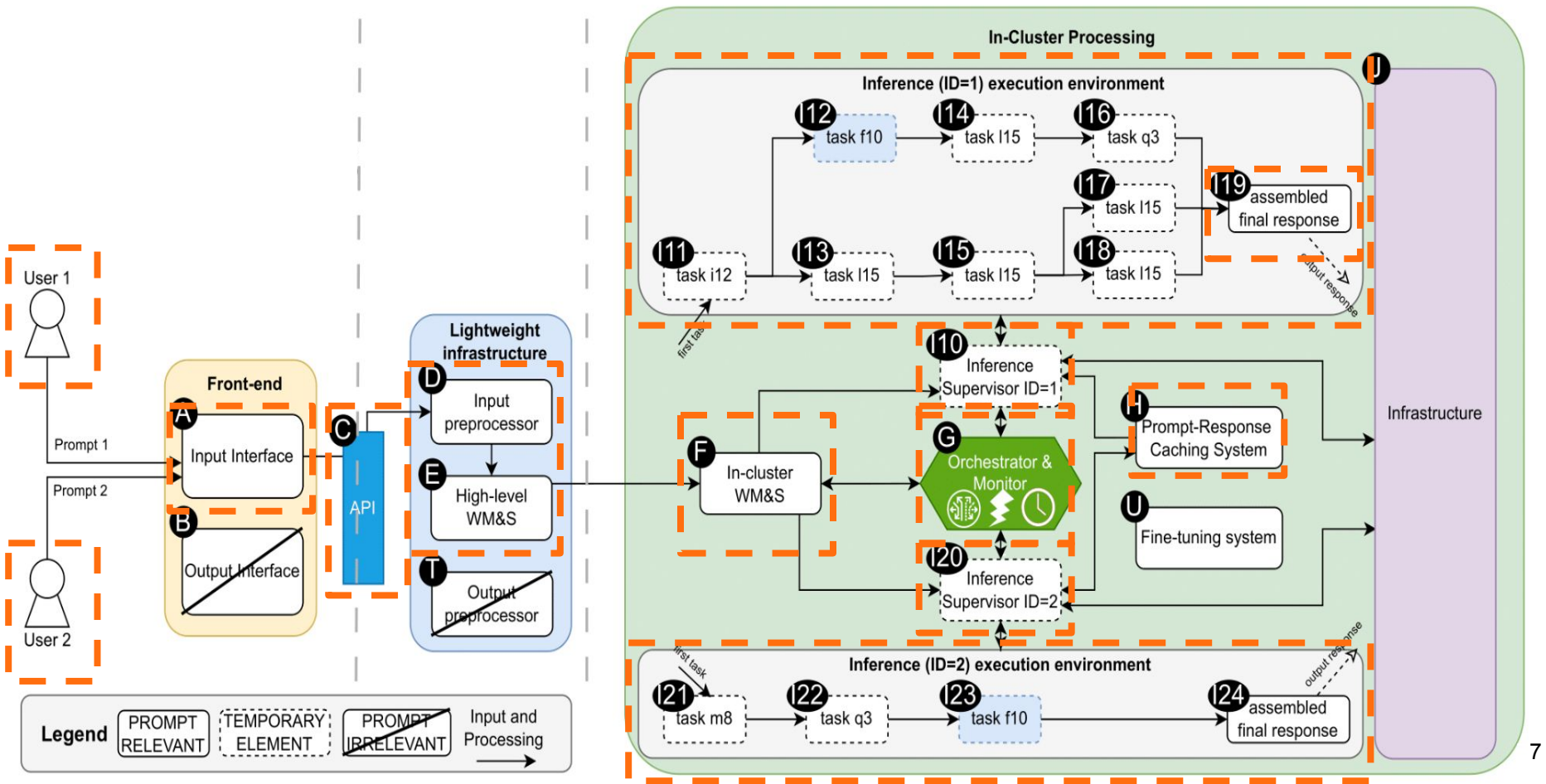
DP5 Inference stages

& more

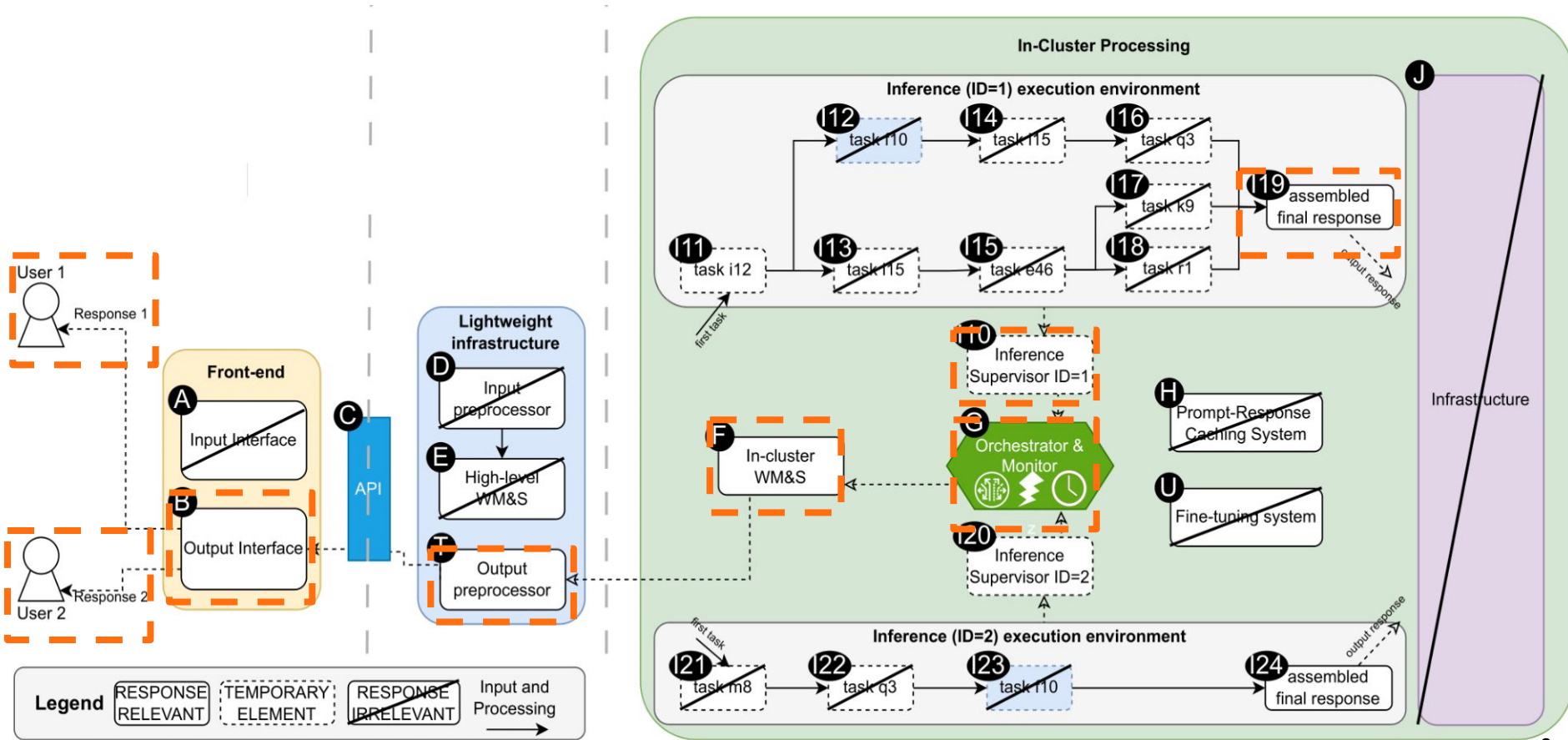
[1] Andreadis et al., A reference architecture for datacenter scheduling, SC'18



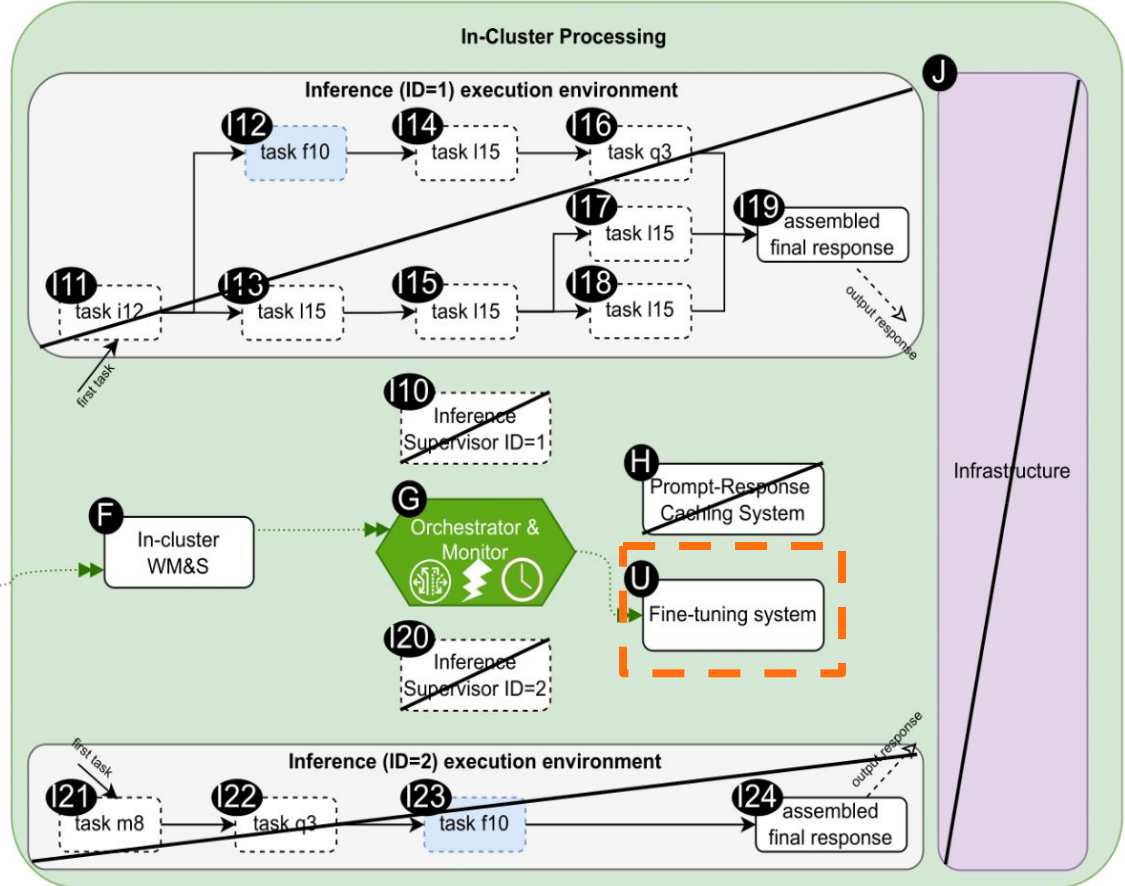
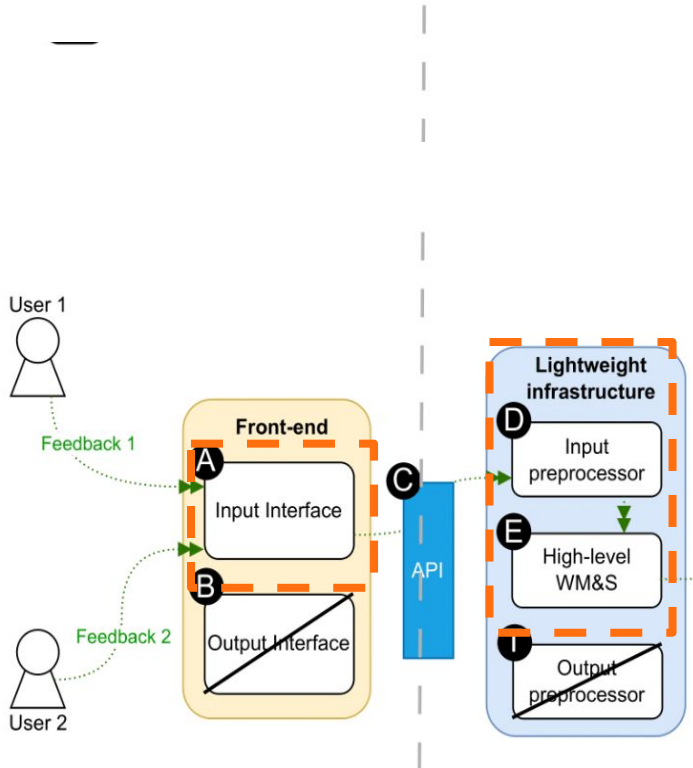
DESIGN - PROMPT WORKFLOW [2/4]



DESIGN - RESPONSE WORKFLOW [3/4]



DESIGN - FEEDBACK WORKFLOW [4/4]

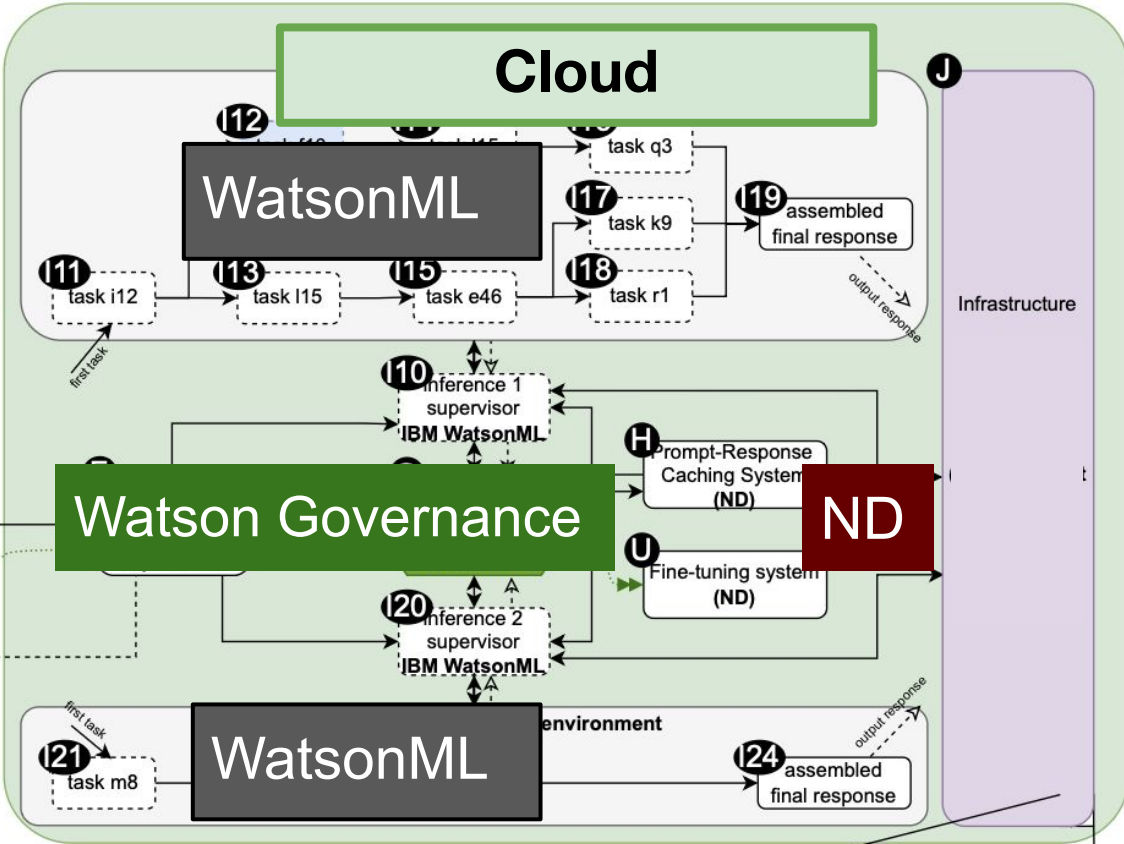
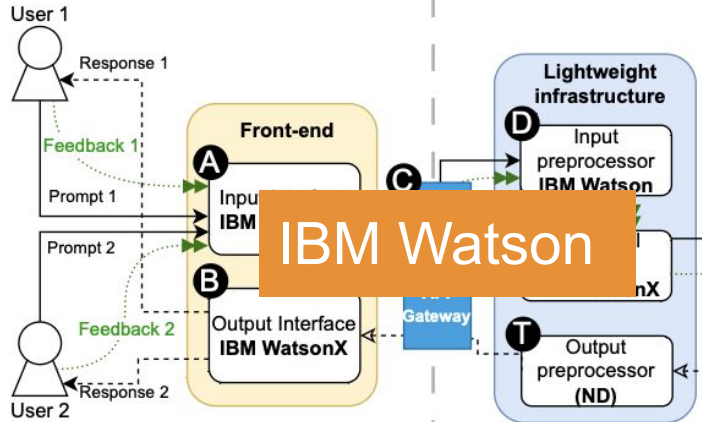


VALIDATION

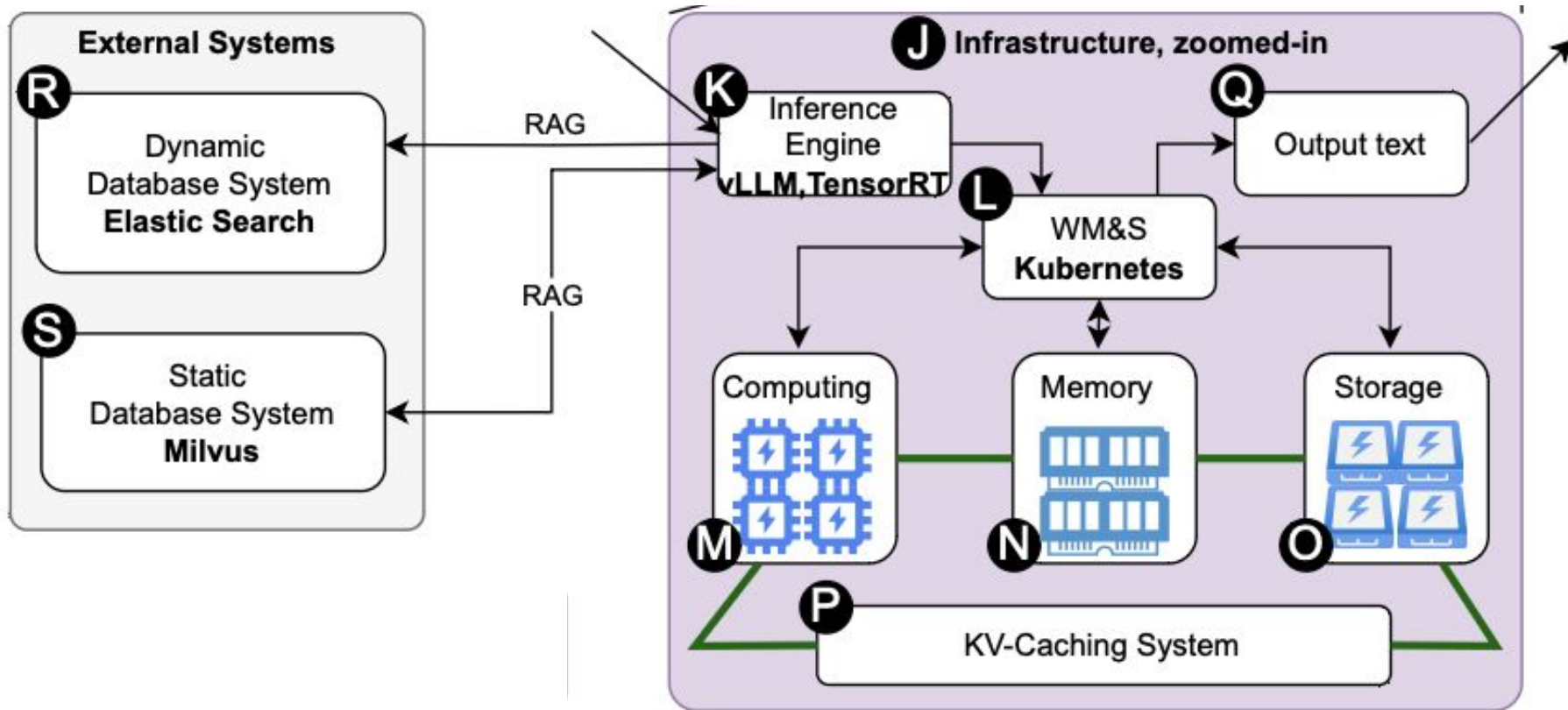
VALIDATION - IBM [1/2]

Endpoint

Edge



VALIDATION - IBM [2/2]



TAKE-HOME MESSAGES (REF. ARCHI.)

SIMULATING LLM ECOSYSTEMS



FROM REFERENCE ARCHITECTURES OF LLM ECOSYSTEMS TO
SIMULATING PERFORMANCE, SUSTAINABILITY, EFFICIENCY

@Large Research
Massivizing Computer Systems

P1

C1: Propose reference architecture
C2: Validate ref. arch. (x3) (-> x4?)

P2

C3: Design LLM inference simulator
C4: Impl. sim, integrate with OpenDC
C5: Trace vLLM deployments
C6: Validate simulator (wip)



Radu Nicolae
@VU Amsterdam
@Large Research



@rnicolae,
mail@radu-nicolae.com

KAVIER:

A SIMULATOR FOR
LLM ECOSYSTEMS UNDER
INFERENCE
(DESIGN, IMPLEMENTATION,
INTEGRATION)

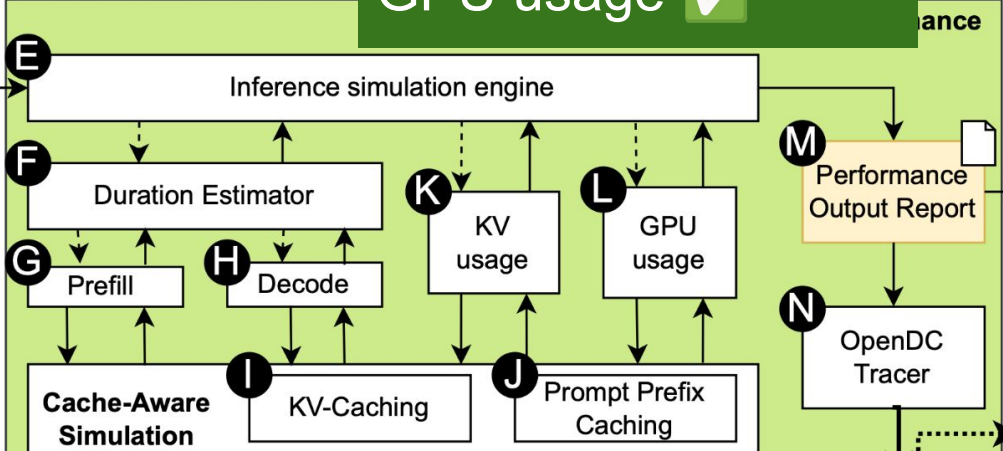
Design of Kavier

Performance Report
GPU usage ✓

Kavier

Kavier Input

- A LLM Configuration
- B GPU Configuration
- C Prompt trace
- D Simulation Setup



Efficiency

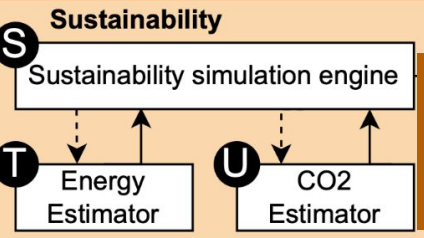
- W Performance-Financial Cost
- X Performance-Sustainability Cost
- Y Efficiency Output Report

Efficiency Report
\$/t/s, CO2/t/s ✓

OpenDC

- O GPU Configuration
- P OpenDC Simulation Setup

- ### OpenDC Input
- Q Tasks
 - R Fragments



Sustainability Report
Wh, CO2 usage ✓



Identical

C4

KAVIER:

IMPLEMENTED A PROTOTYPE,
INTEGRATED WITH OPENDC,
OPEN-SOURCE & OPEN-SCIENCE

TRACER

LLM Ecosystem Tracer

Challenge: no traces decode/prefill time per decode/prefill length

Approach: accessed clusters, deployed vLLM, measured performance



SURF

vLLM

A10

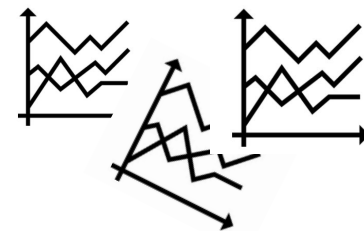


vLLM

A4000,
A6000,
A100 ++

```

  ✓ folder tracer
    > folder _data
    > folder clients
    > folder keys
    ≡ .env
    🐍 __init__.py
    🐍 config.py
    🐍 logger.py
    🐍 main.py
    🐍 prompt.py
    M↓ README.md
    🐍 sampler.py
  
```



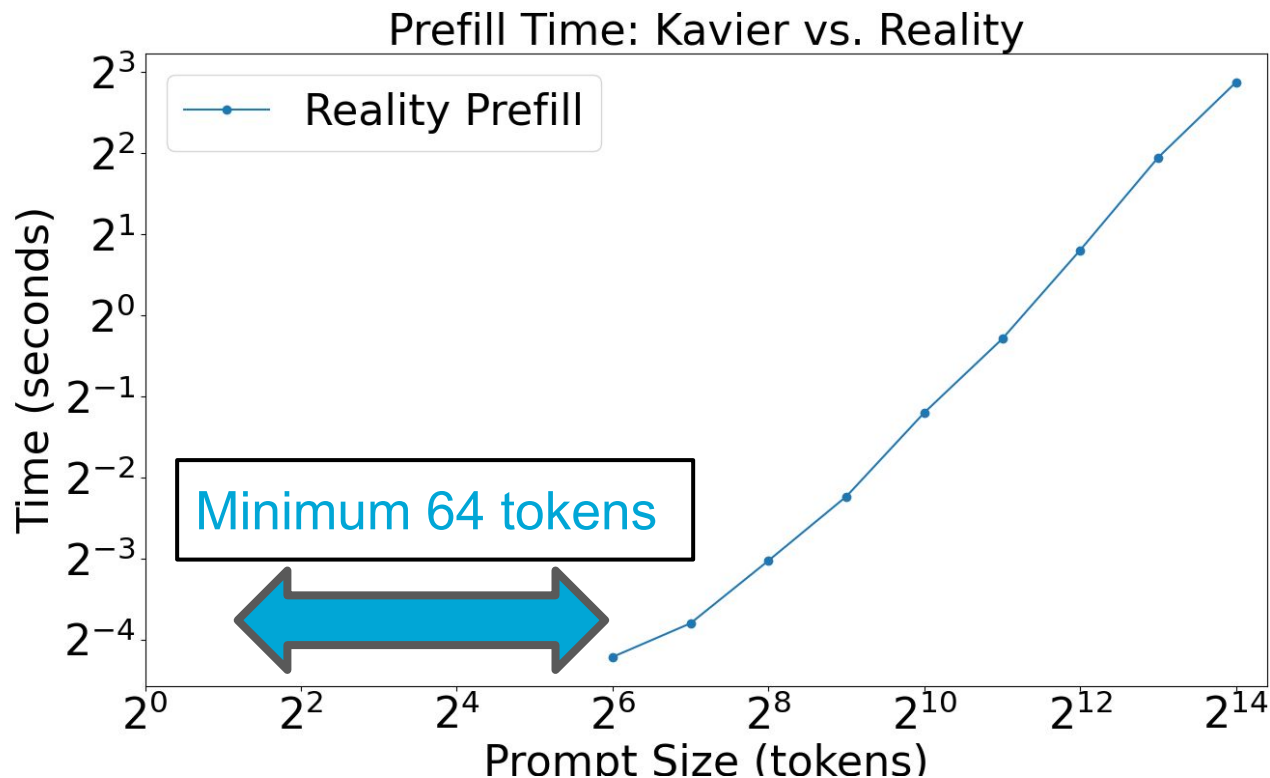
Traces released
open-source &
open-science

VALIDATION

Validation - Prefill [1/4]

Experiment setup:

- SURF
- A10
- vLLM



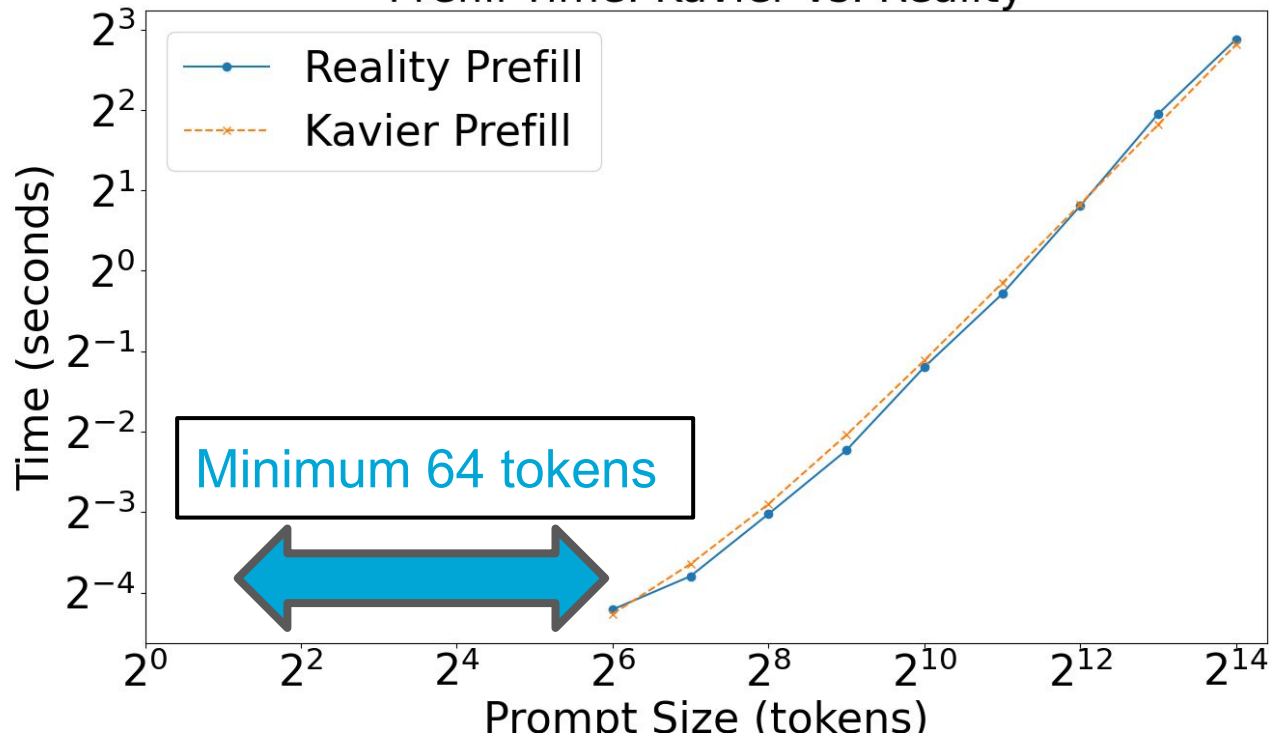
Validation - Prefill [2/4]

MAPE: 7.58%

Prefill Time: Kavier vs. Reality

Experiment setup:

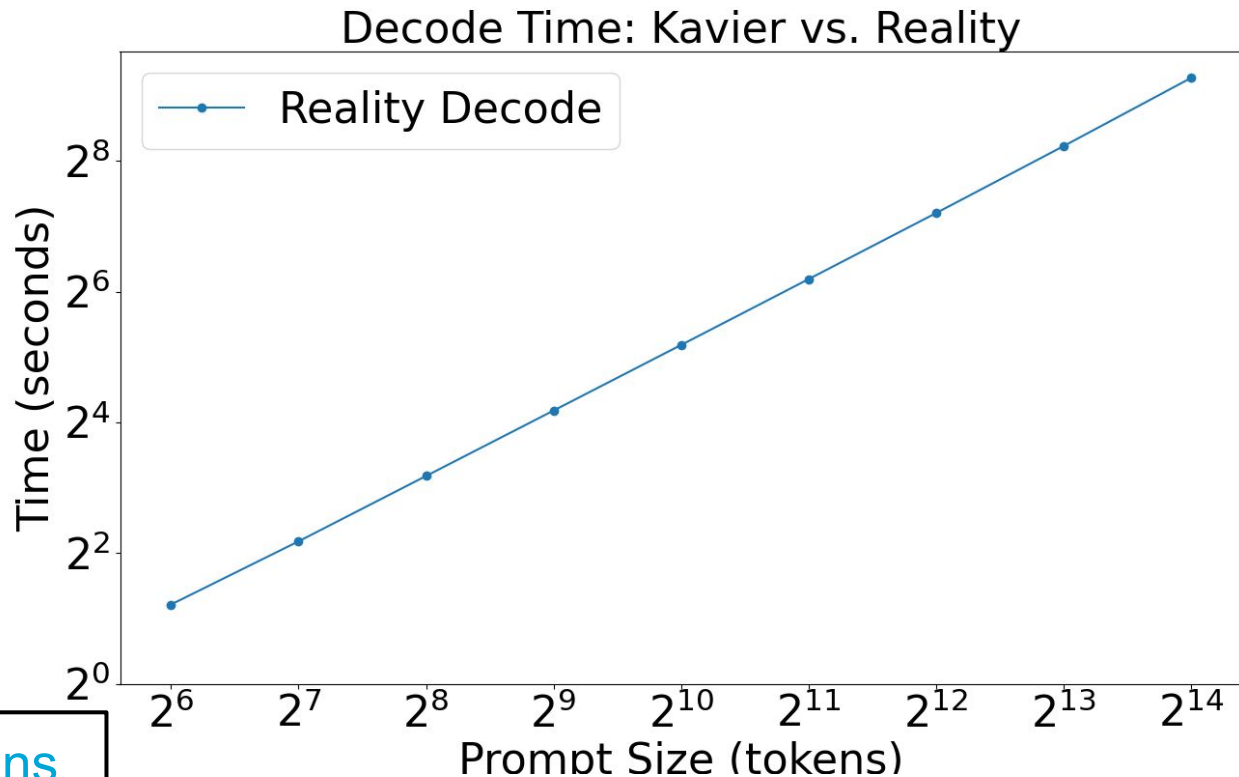
- SURF
- A10
- vLLM



Validation - Decode [3/4]

Experiment setup:

- SURF
- A10
- vLLM



Starting at 64 tokens

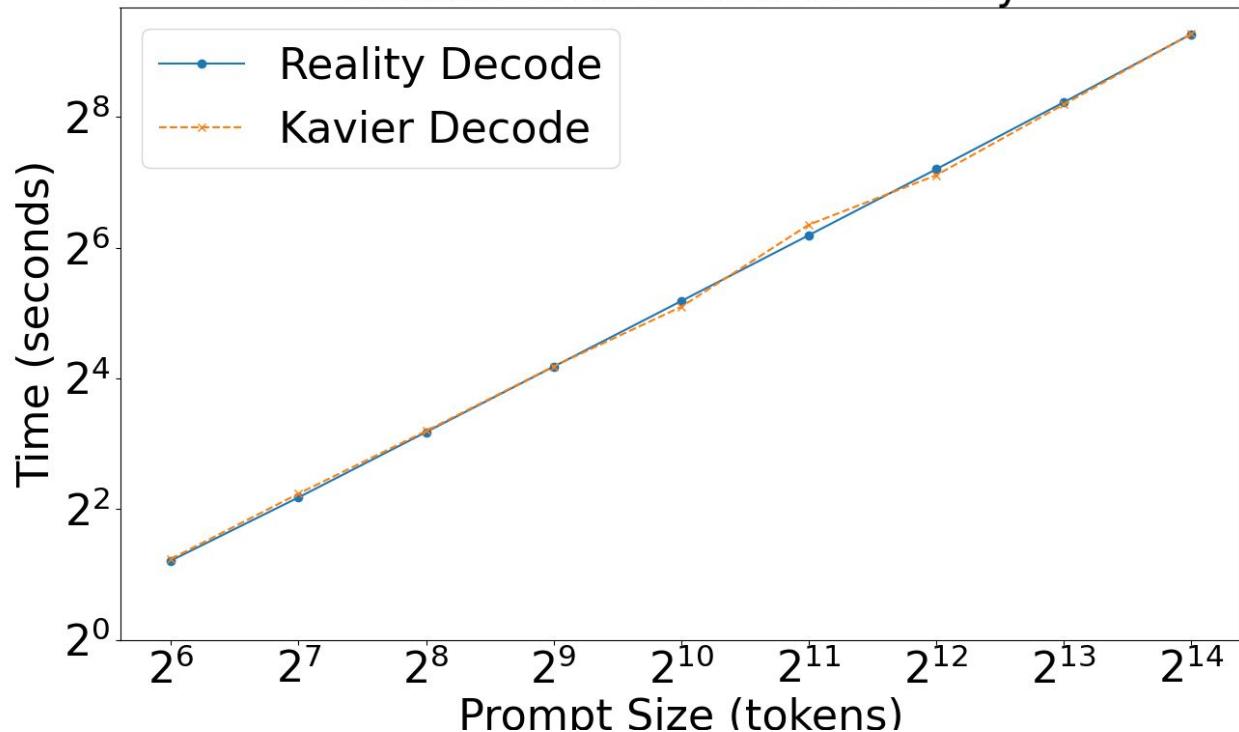
Validation - Decode [4/4]

MAPE: 3.75%

Experiment setup:

- SURF
- A10
- vLLM

Decode Time: Kavier vs. Reality



SIMULATING LLM ECOSYSTEMS



FROM REFERENCE ARCHITECTURES OF LLM ECOSYSTEMS TO
SIMULATING PERFORMANCE, SUSTAINABILITY, EFFICIENCY


@Large Research
Massivizing Computer Systems

P1

C1: Propose reference architecture
C2: Validate ref. arch. (x3)



Radu Nicolae
@VU Amsterdam
@Large Research

 @rnicolae,
mail@radu-nicolae.com

P2

C3: Design LLM inference simulator
C4: Impl. sim, integrate with OpenDC
C5: Trace vLLM deployments
C6: Validate simulator (wip)