Trust arises from **knowledge of origin** as well as from **knowledge of functional capacity**.

*Trustworthiness - Working Definition* by *David G. Hays, 1979*

Task

Expression

Expectation
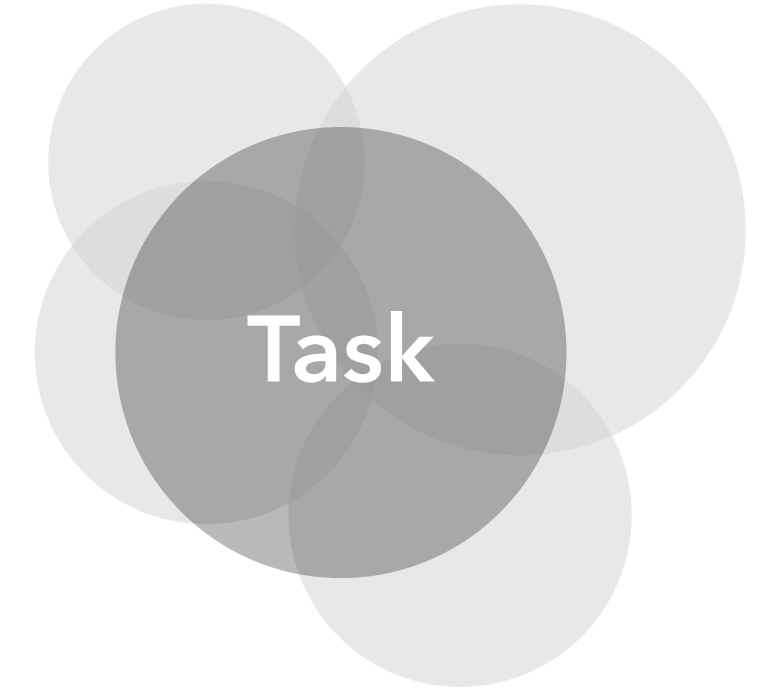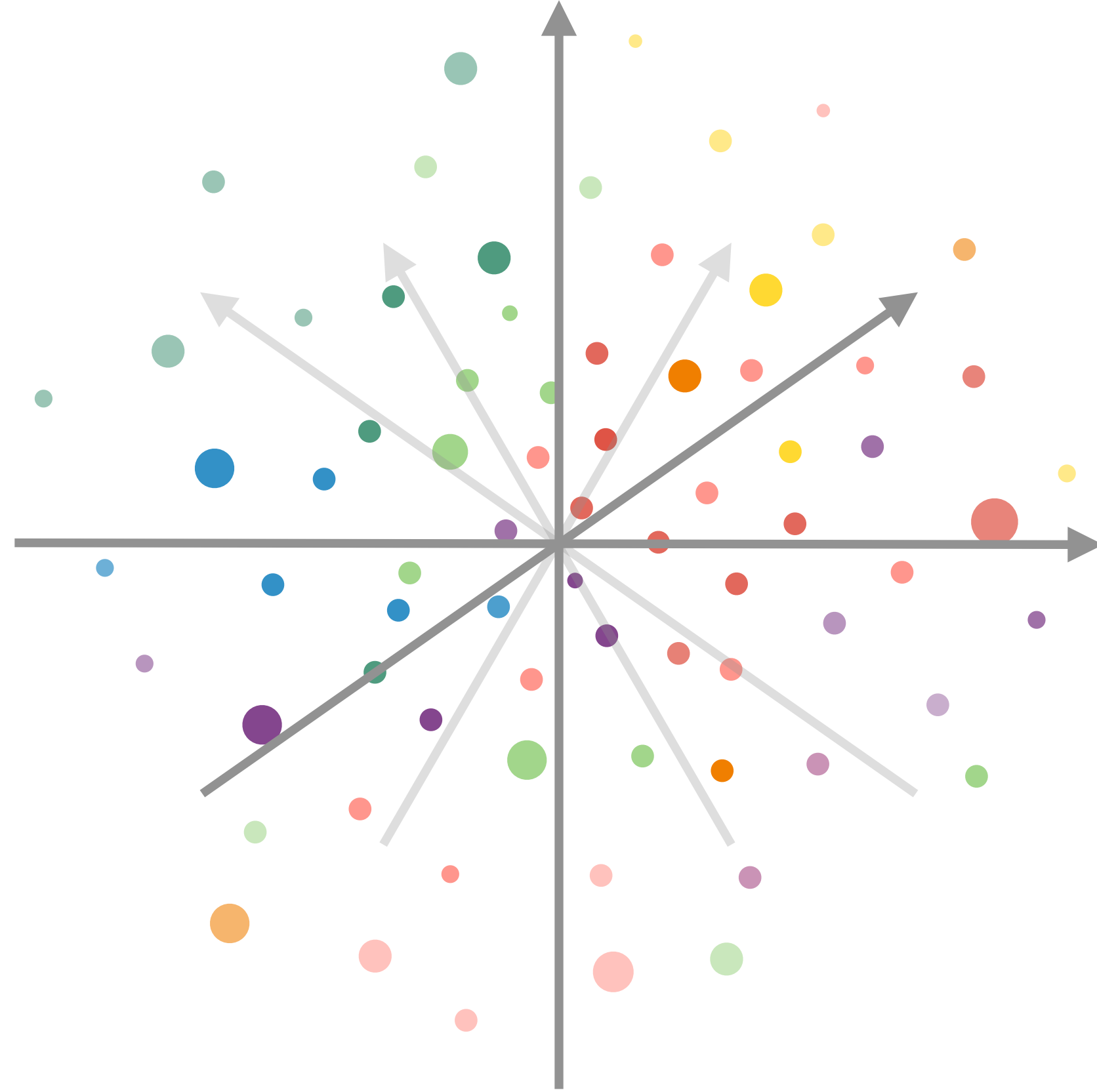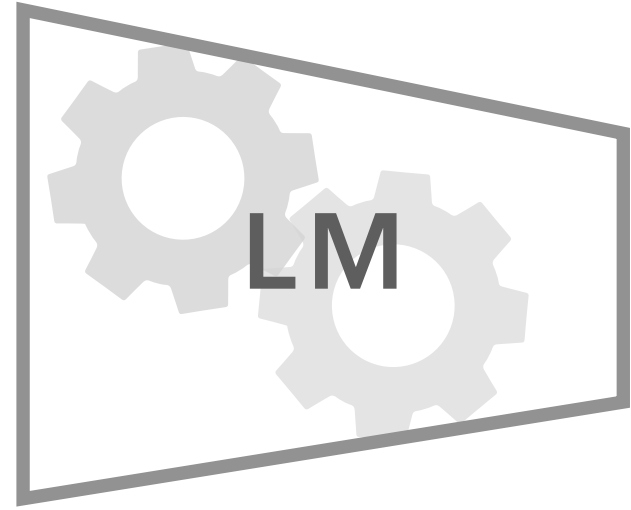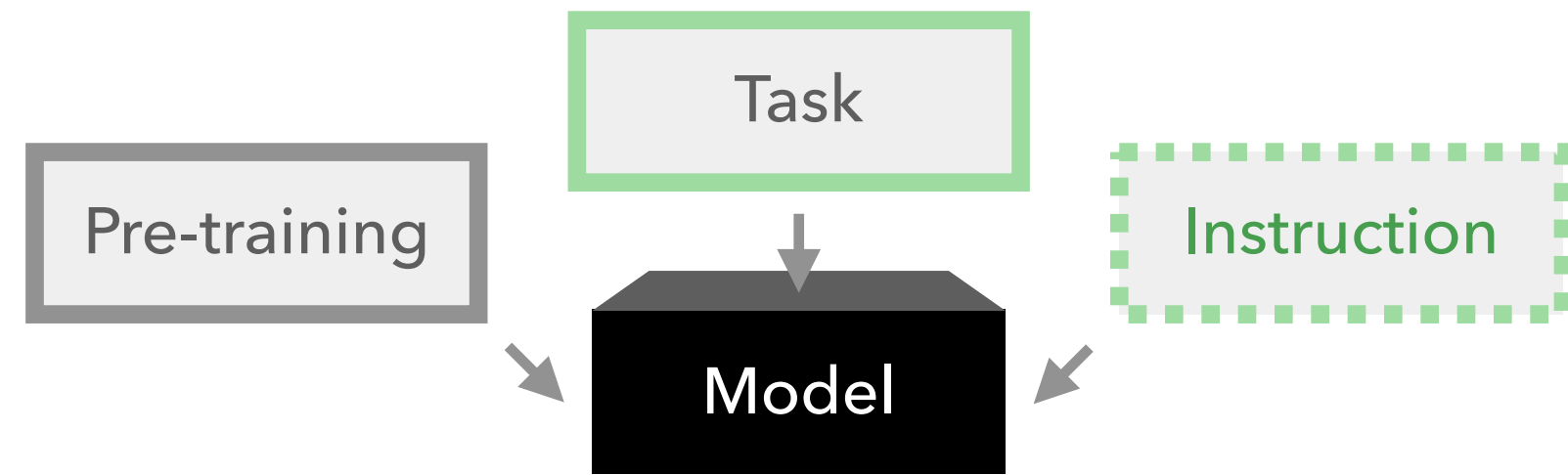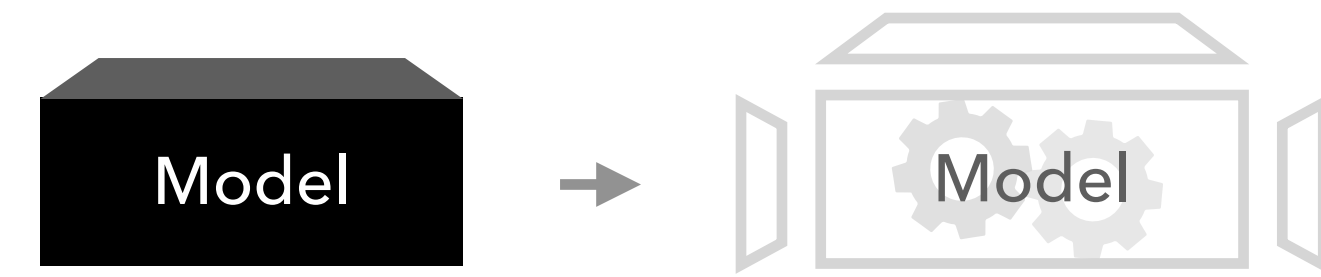
Datasets

D1 Knowledge about Model Input

D2 Knowledge about Model Behaviour

D3 Knowledge of Evaluation Protocols

D4 Knowledge of Data Origin

D1 Knowledge about Model Input

Pre-training

Task

Instruction

Model

Output

+ Correct, Truthful, Up-to-Date
− Hallucinated, Biased, Outdated

Feature Engineering
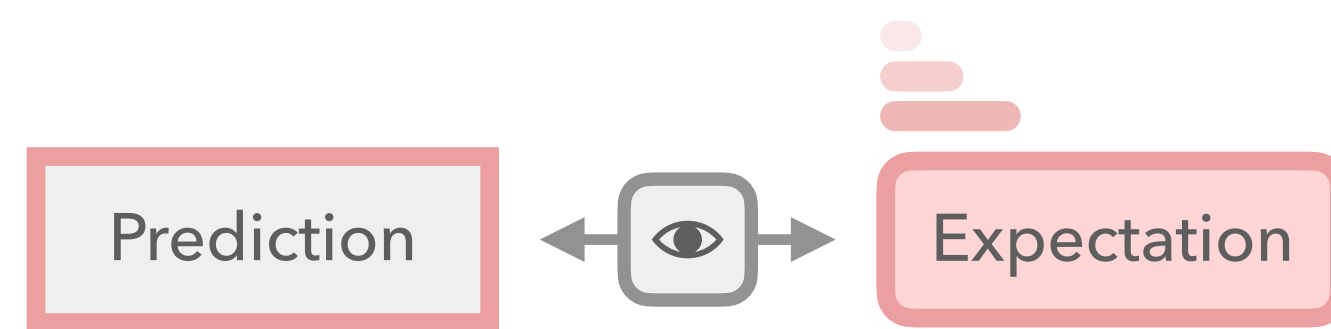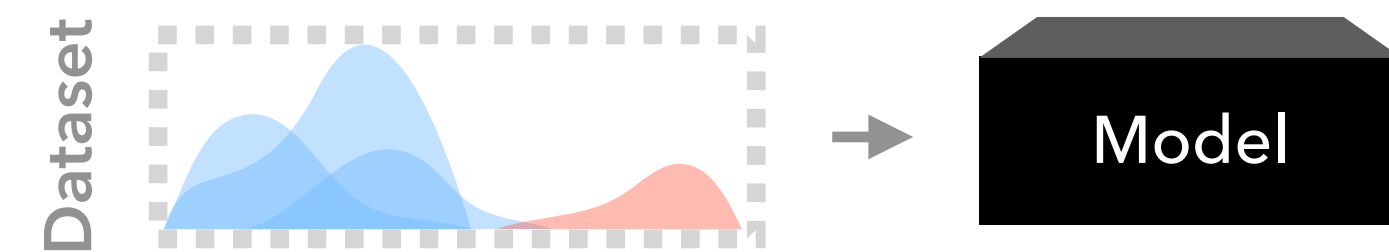
1.

Representation Learning

2.

Input and External Knowledge

Task Instructions

Task Examples

## D1 Knowledge about Model Input

Pre-training

Task

Instruction

Model

Which internal and external knowledge (ingredients) is used to generate output?

## Output

+ Correct, Truthful, Up-to-Date
− Hallucinated, Biased, Outdated

## Internal Knowledge

World Knowledge

Common Sense

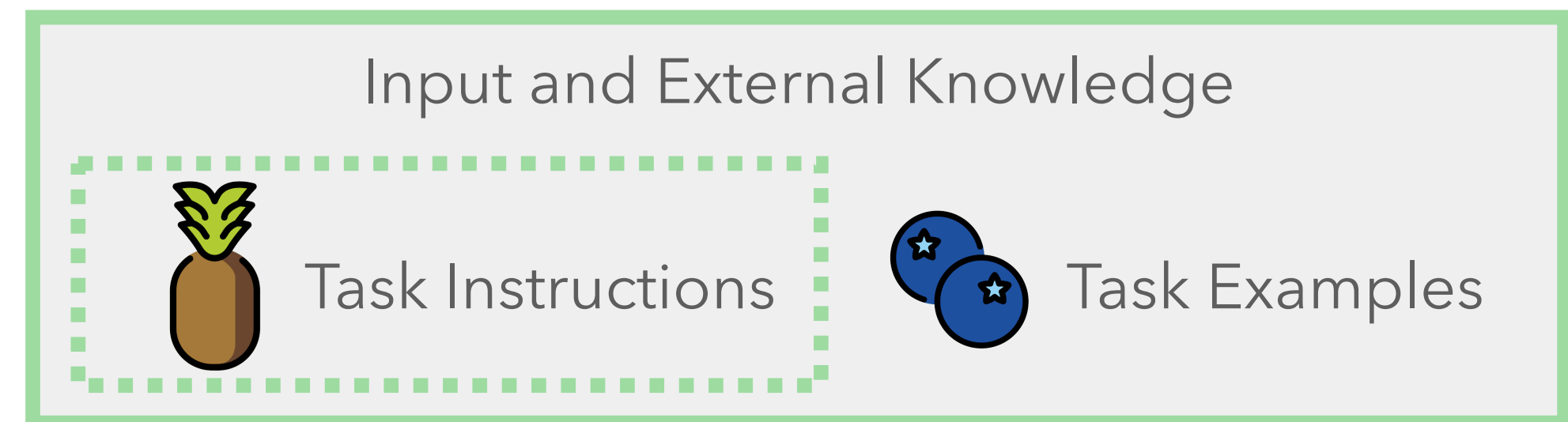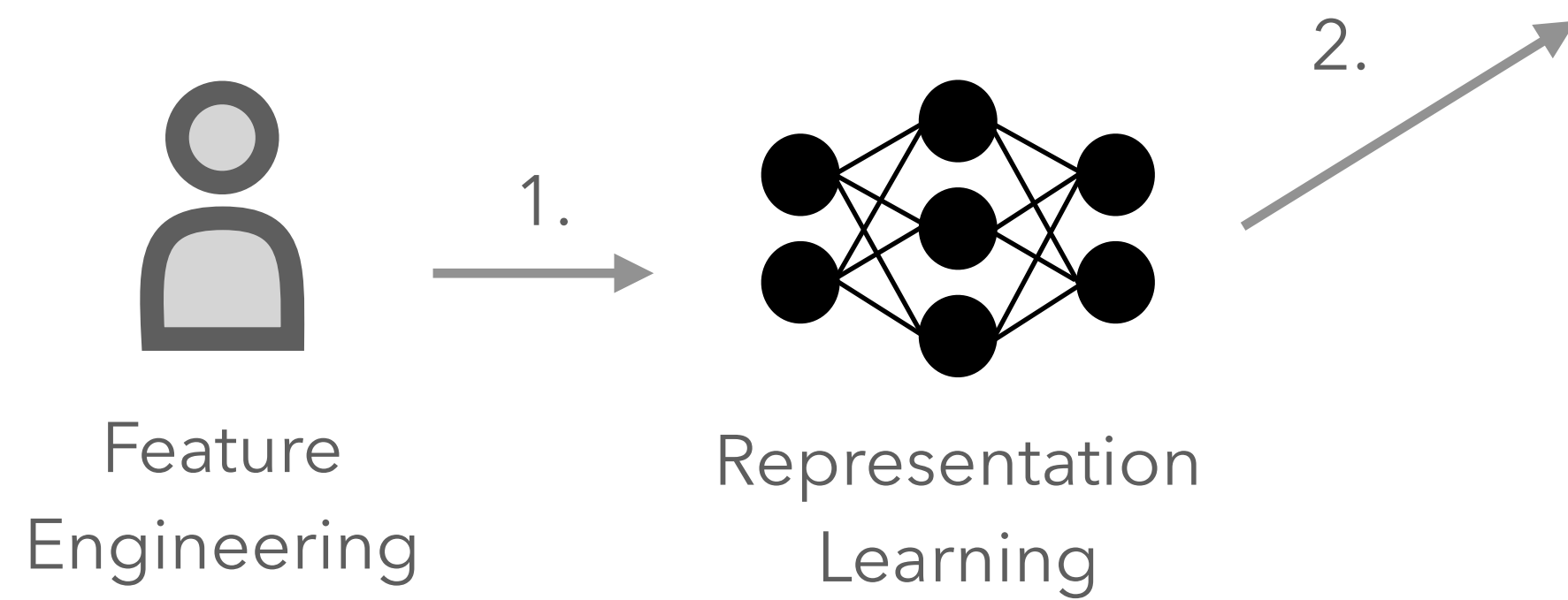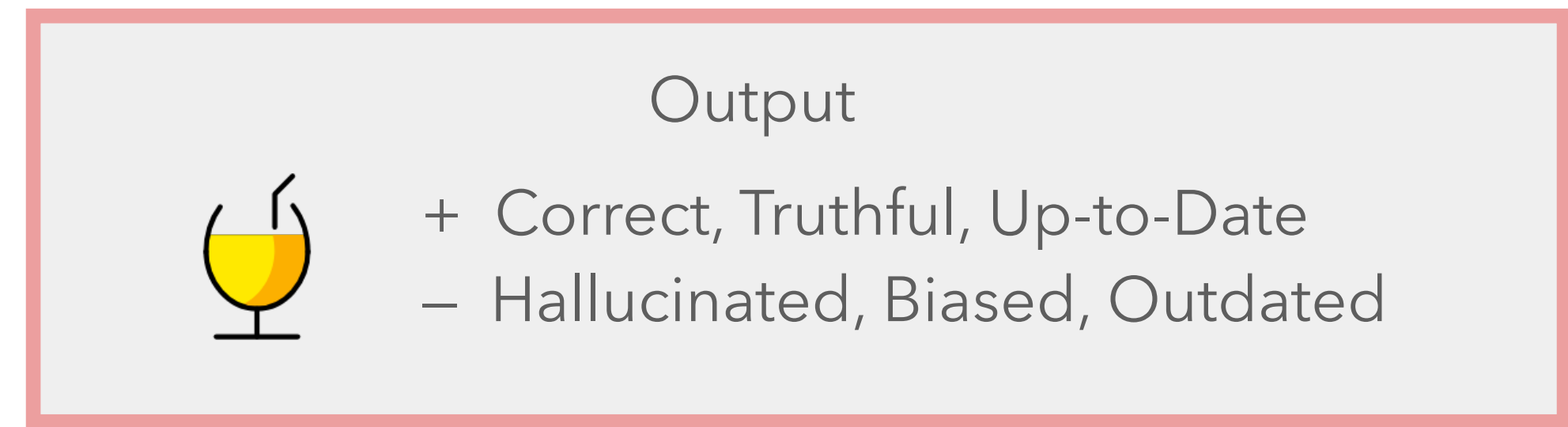Propaganda

Facts

## Input and External Knowledge

Task Instructions

Task Examples

D2 Knowledge about Model Behaviour

Model → Model

What skills are employed to process ingredients into model output?

Output

+ Correct, Truthful, Up-to-Date
− Hallucinated, Biased, Outdated

**Internal Knowledge**

World Knowledge

Common Sense

Propaganda

Facts

**Model Skills**

Linguistic and Logical Reasoning

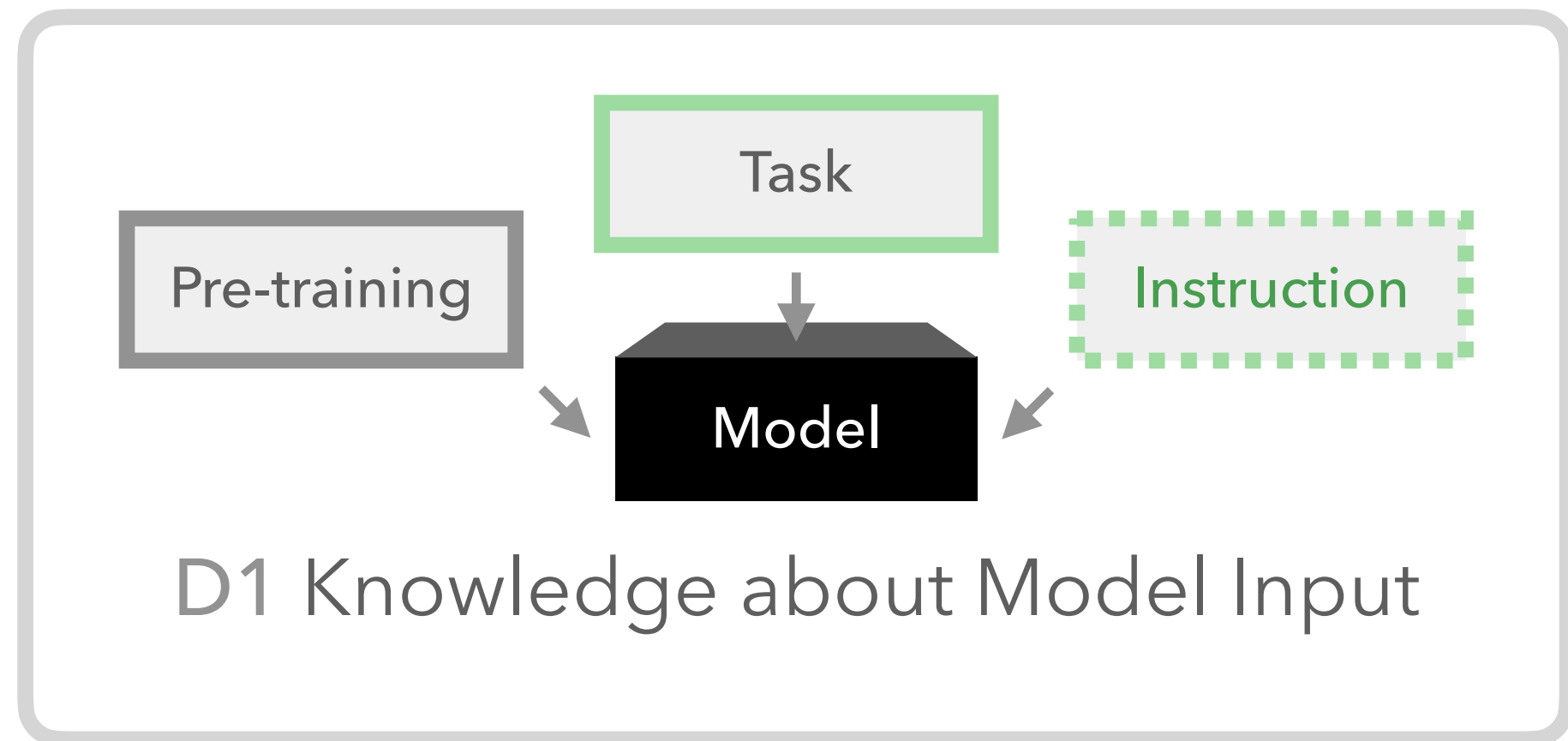Shortcuts

Memorization

Input and External Knowledge

Task Instructions

Task Examples

**D3 Knowledge of Evaluation Protocols**

Prediction ↔ 👁 ↔ Expectation
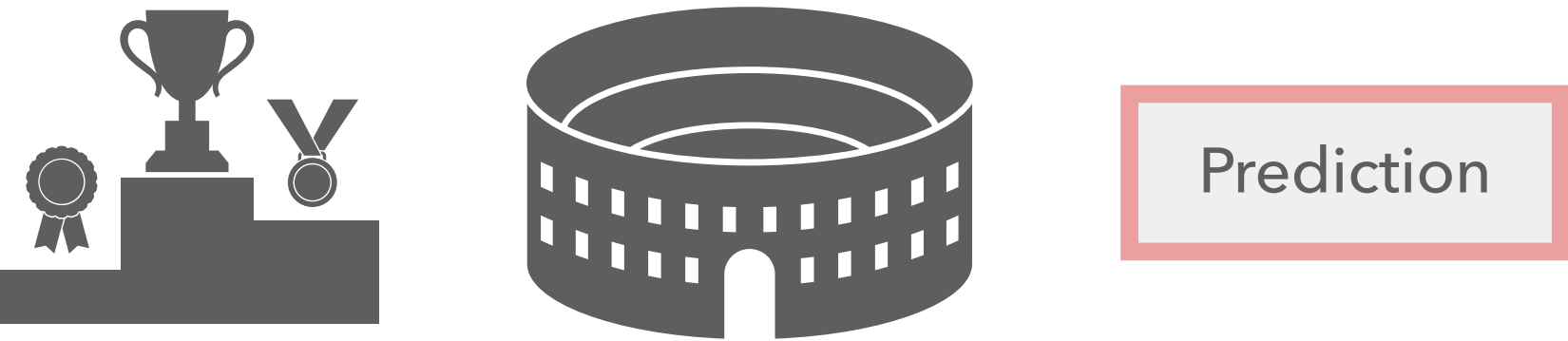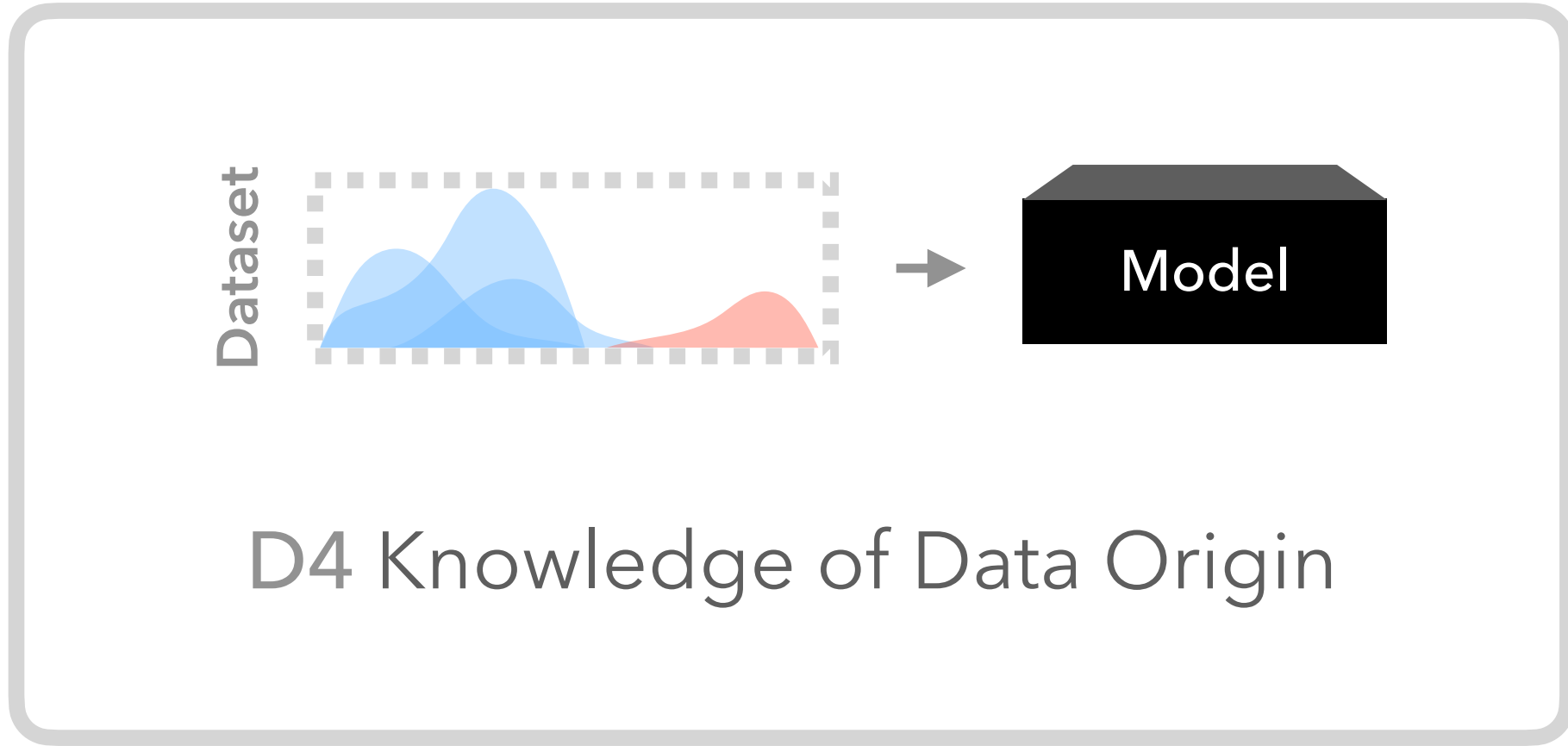
**Baselines:** Super-human Performance?

LLM ⇢ 👤⚖ Expectation
⇢ 👤➕
⇢ 👤 Crowdsourcing Workers

**Benchmarking:** Leaderboard-style?

🏆 🏟 Prediction

Dataset

Model

D4 Knowledge of Data Origin

Compartmentalized NLP

Train | Dev | Test

Hate Speech, Biases, Copyright?

Test

LLM

Current Trend in NLP

LLM

Pre-training + RLHF

# Conclusion

Trust arises from **knowledge of origin** as well as from **knowledge of functional capacity**.
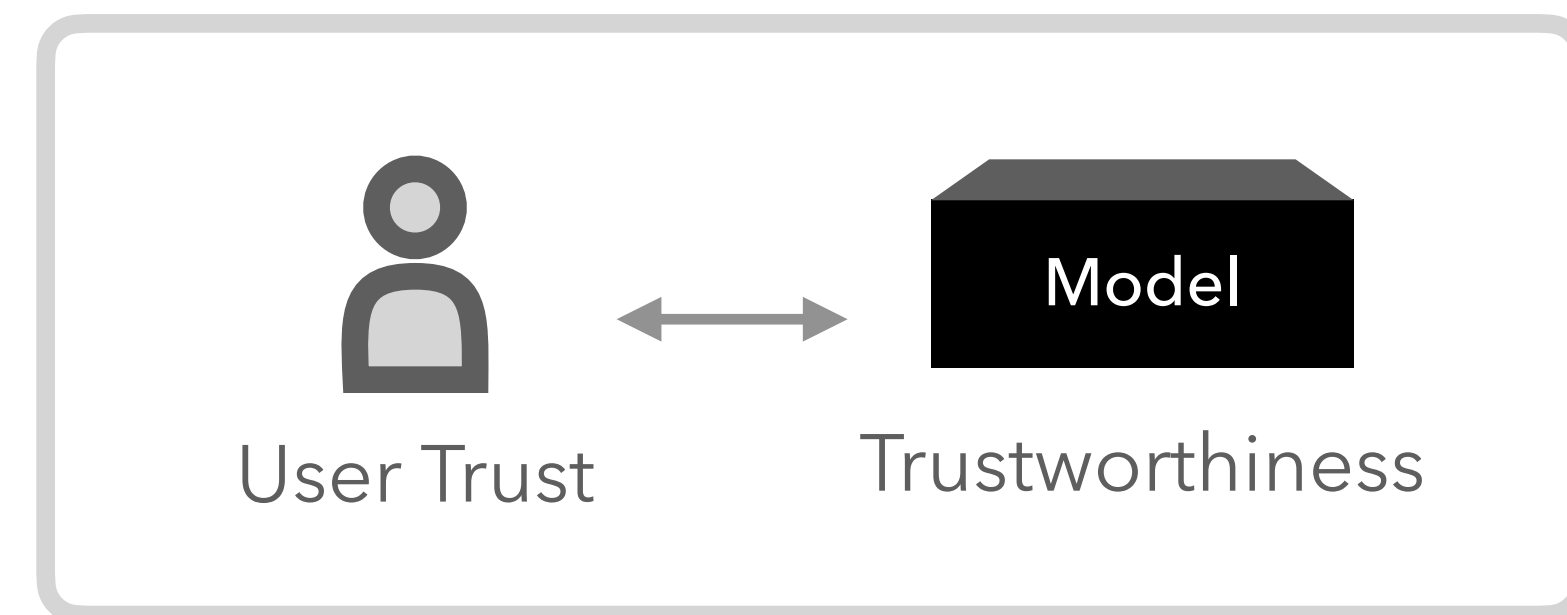*Trustworthiness - Working Definition* by *David G. Hays, 1979*

**D1** Knowledge about Model Input

**D2** Knowledge about Model Behaviour

**D3** Knowledge of Evaluation Protocols

**D4** Knowledge of Data Origin

User Trust          Model          Trustworthiness
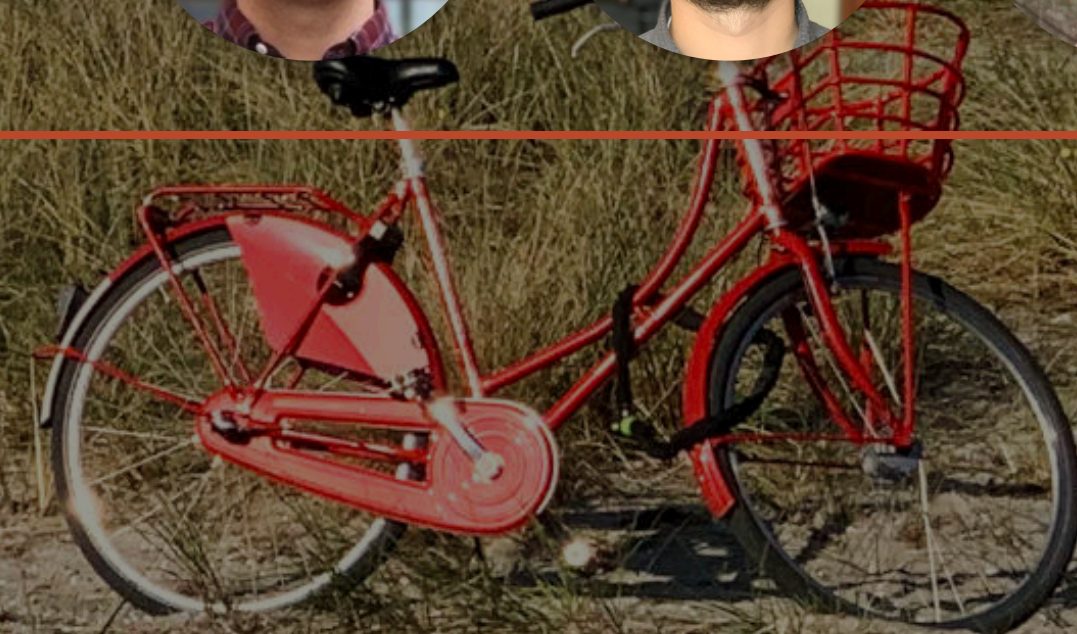
What Can We Do to Gain Trust Now and in Future? (See paper)

- Explain Skills Required versus Skills Employed.

- Facilitate Representative and Comparable Qualitative Analysis.

- Be Explicit about Data Provenance.

**Thanks**

See you in Singapore!

EMNLP

IT-UNIVERSITETET I KØBENHAVN