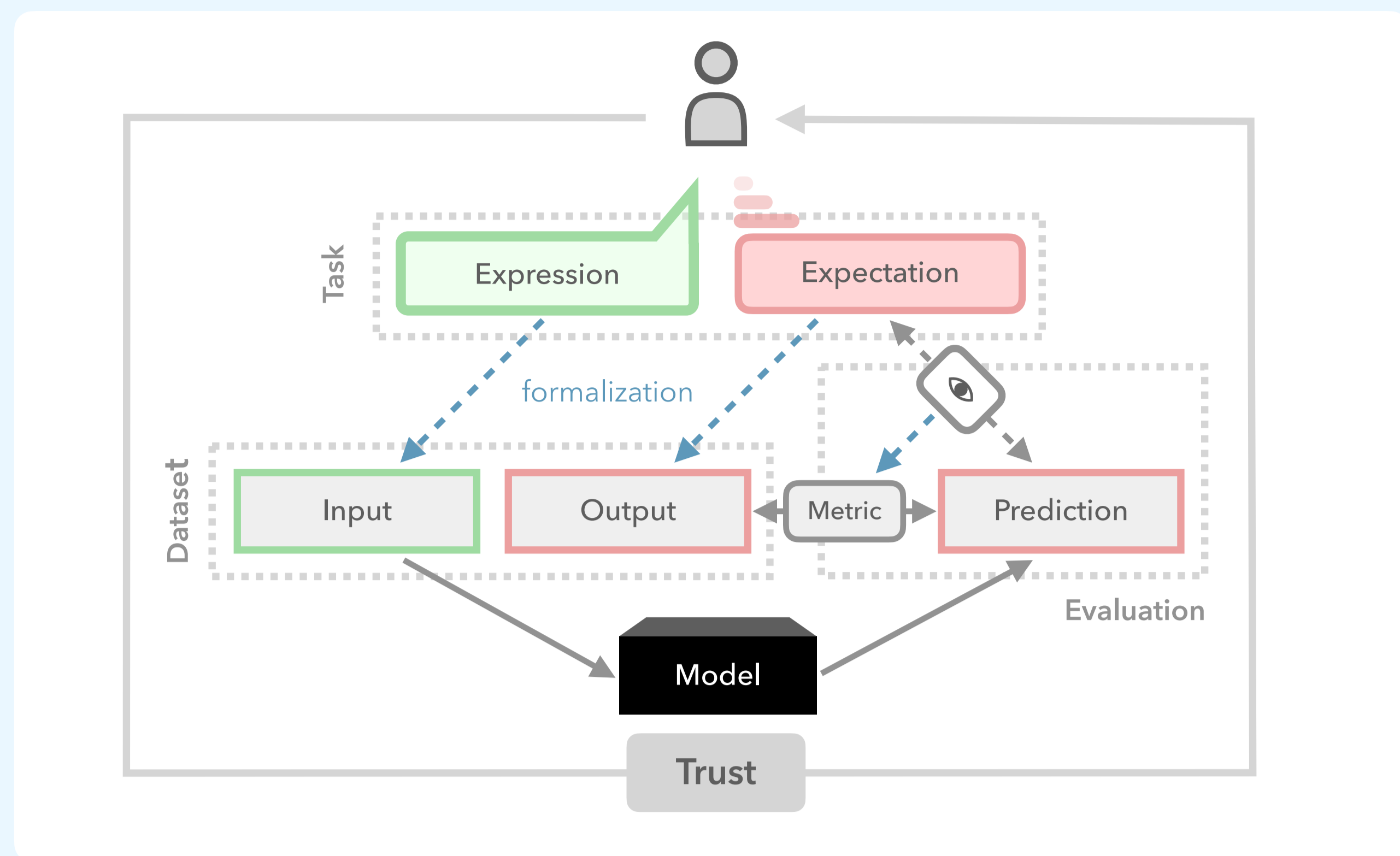


# Establishing Trustworthiness: Rethinking Tasks and Model Evaluation

Robert Litschko\*, Max Müller-Eberstein\*, Rob van der Goot, Leon Weber, Barbara Plank

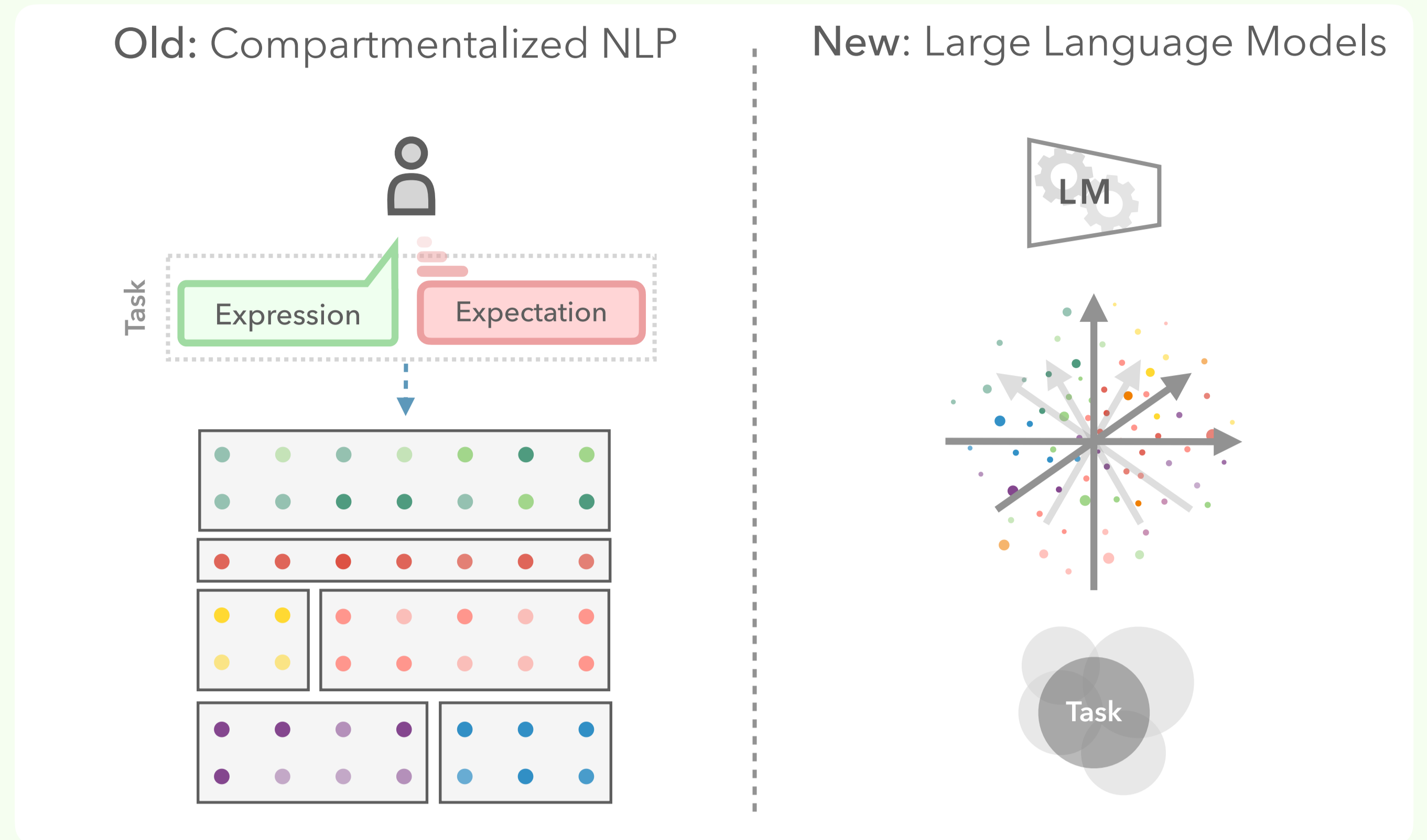
## Working Definition



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

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

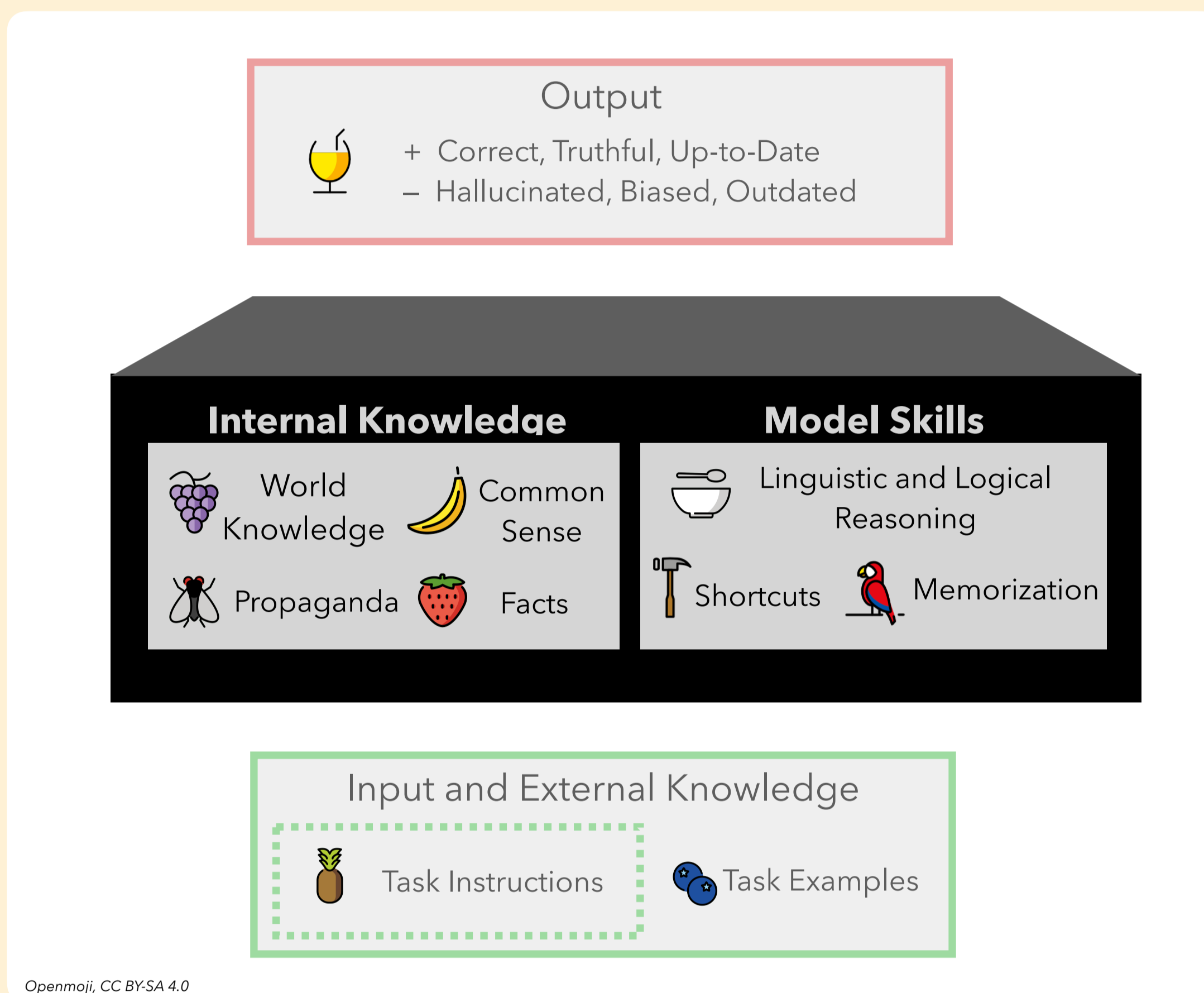
## Paradigm Shift



Task-specific Models and Evaluation Protocols

Massive Multi-task and Multilingual Learning

## Trustworthiness Desiderata



Openmoji, CC BY-SA 4.0

D1 Knowledge about Model Input

- Over time, we have lost the ability to **inspect intermediate outputs** and **control model inputs**.
- What "ingredients" (🍇, 🍌, 🍓, 🐛) are used to generate output?

D2 Knowledge about Model Behaviour

- Over time, we lost the ability to **interpret decision boundaries** and **model behaviour**.
- Which skills (🍷, 🛠️, 🐦) are employed to process ingredients into outputs?

D3 Knowledge of Evaluation Protocols

- Today, LLMs are used to solve **tasks outside of the benchmark**: user-formulated tasks via instructions.
- What do human-level comparisons and leaderboards tell us about **model capabilities**?

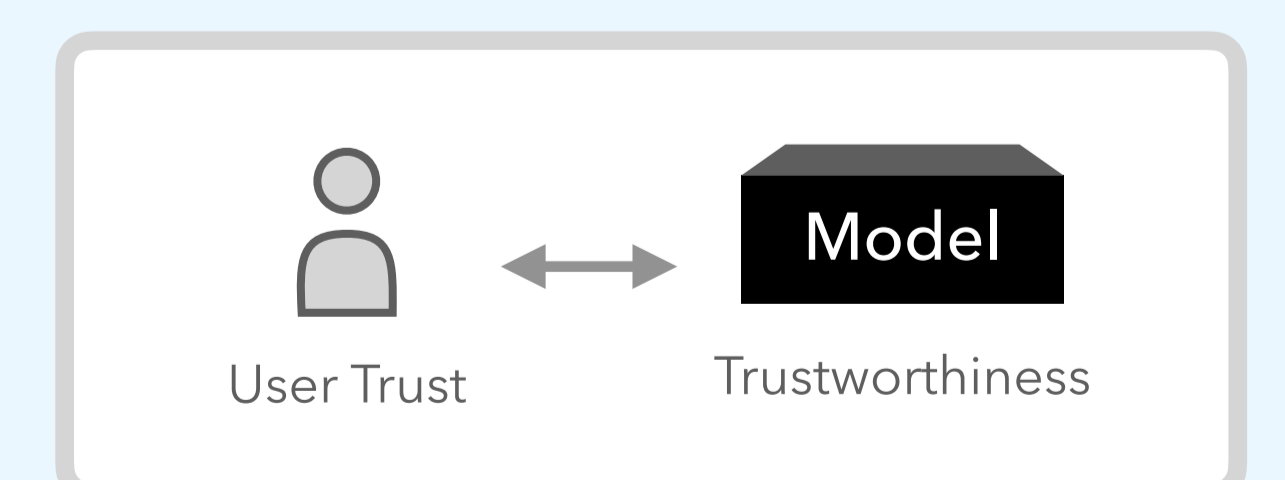
D4 Knowledge of Data Origin

- Is **model performance reflective of its capabilities** or overestimated due to data leakage?
- Additional risks** from unknown provenance include hate speech, biases, copyright violations, ...

## What Can We Do to Gain Trust?

<p>Explain <b>skills required</b> (🍷) vs. <b>skills employed</b> (🛠️, 🐦)</p> <ul style="list-style-type: none"> <li><b>Linguistically-motivated:</b> <ul style="list-style-type: none"> <li>Probing tasks, Checklists</li> <li>Linguistic Profiling</li> </ul> </li> <li><b>Model-based:</b> Attribute skills to parameter regions.</li> <li><b>Interpretability methods.</b></li> </ul>	<p>Facilitate <b>Representative and Comparable Qualitative Analysis</b></p> <ul style="list-style-type: none"> <li><b>Faceted quantitative analysis.</b></li> <li><b>Standardized qualitative evaluation protocols.</b></li> <li><b>Expert-curated diagnostics annotations</b> w.r.t cognitive abilities required to solve task.</li> </ul>	<p>Be explicit about <b>data provenance</b>.</p> <ul style="list-style-type: none"> <li>Opt for <b>cross-X evaluation</b>.</li> <li><b>Closed-source models</b> (e.g ChatGPT) typically evolve over time and have unknown data provenance.</li> <li>→ <b>Untrustworthy baselines</b></li> </ul>
--	---	---

## User Trust



- Trustworthiness:** Knowledge about LLM's functional capacity and origin.
- User Trust:** What do users do with the model output across multiple interactions (e.g., verify, fact check, revise, accept)?