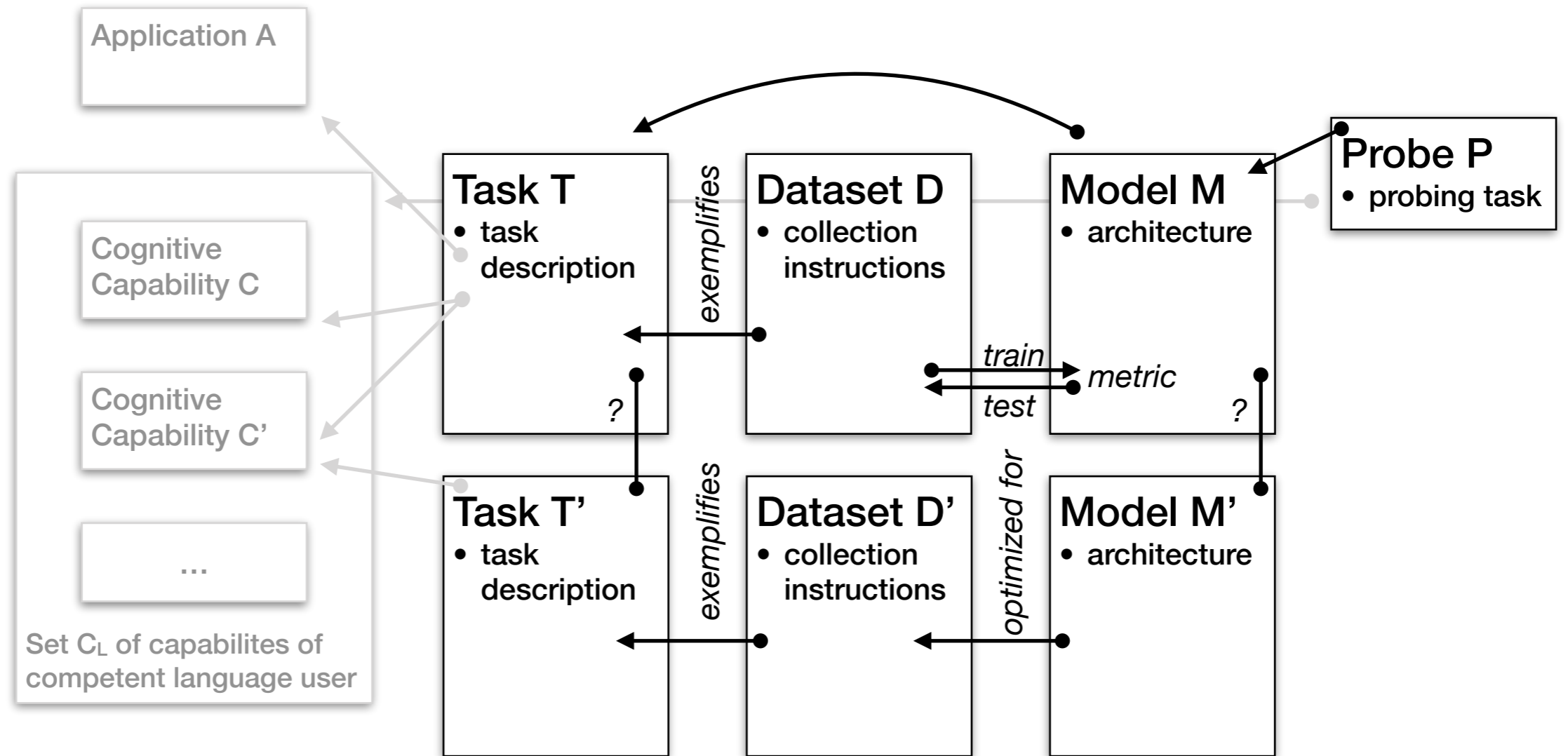
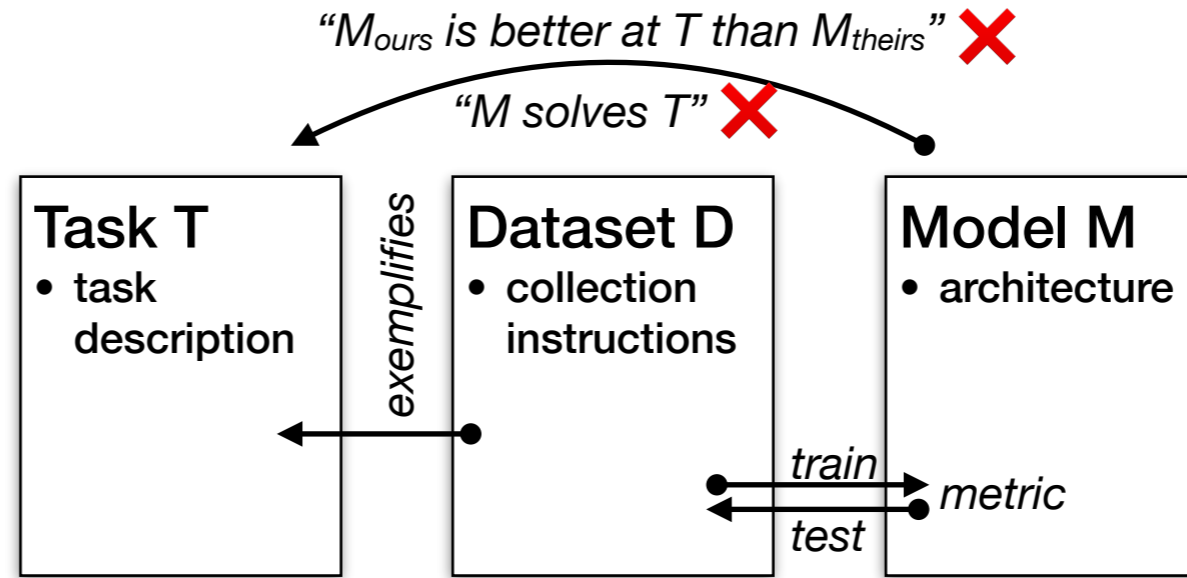


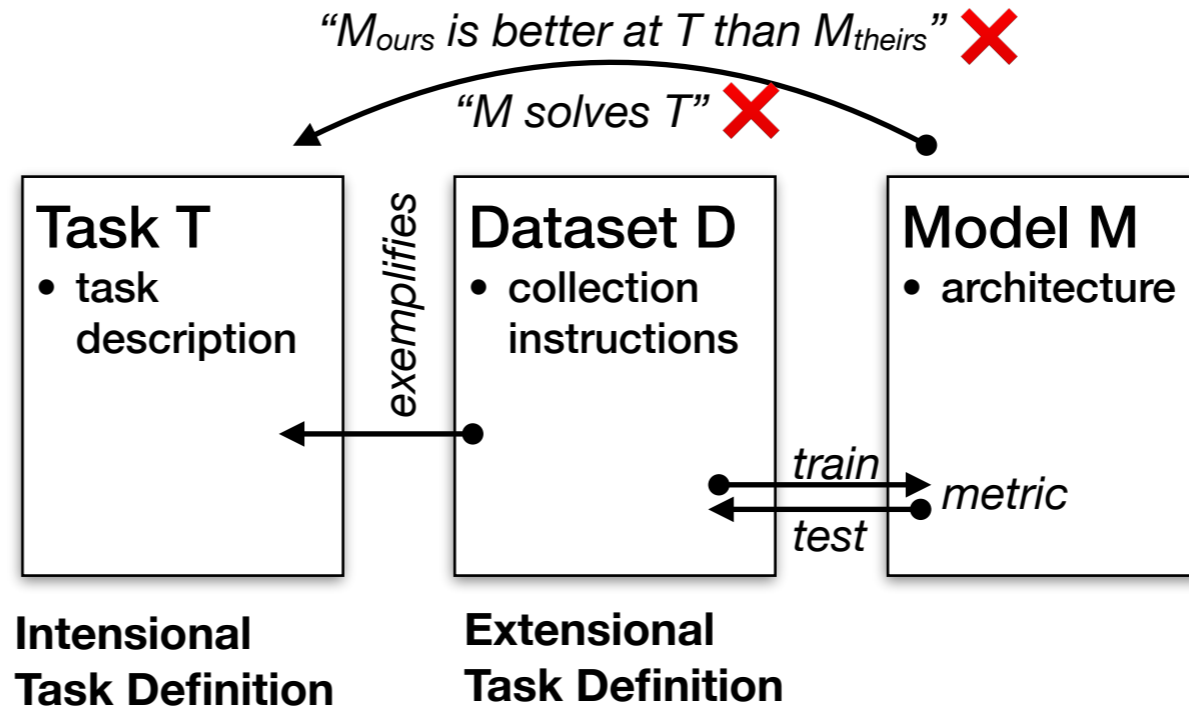
# Targeting the Benchmark

## On Methodology in Current NLP Research

David Schlangen - Dept. Linguistics - University of Potsdam, Germany - [david.schlangen@uni-potsdam.de](mailto:david.schlangen@uni-potsdam.de)

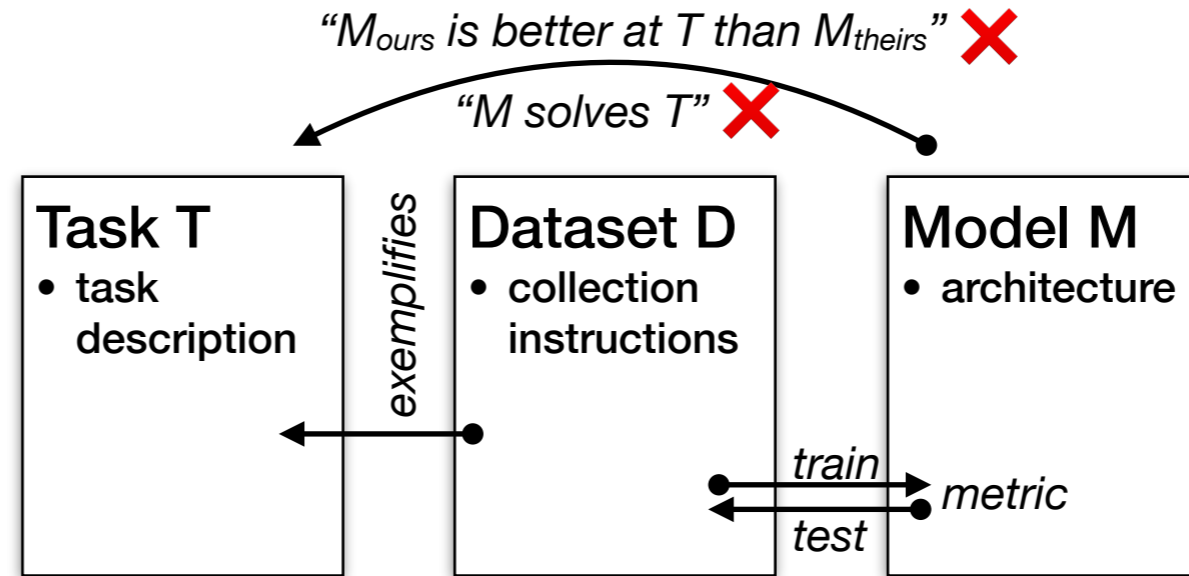






*Image captioning:  
given image (of  
any kind), return  
textual descr. of  
its content.*

$(x, y)$



**Intensional Task Definition**

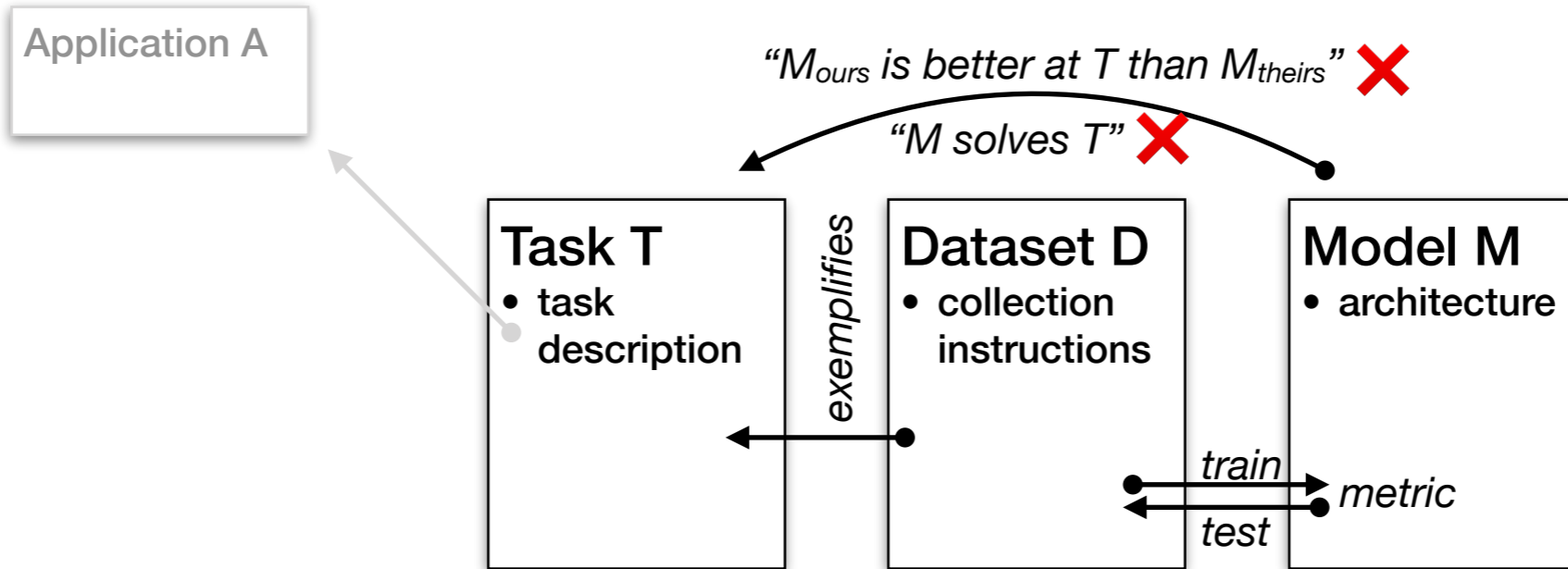
*Image captioning:  
given image (of any kind), return textual descr. of its content.*

**Extensional Task Definition**

$(x, y)$

**Quality Control**

- *internal verification*
- *validation (how general? exploitable?)*



**Intensional Task Definition**

*Image captioning: given image (of any kind), return textual descr. of its content.*

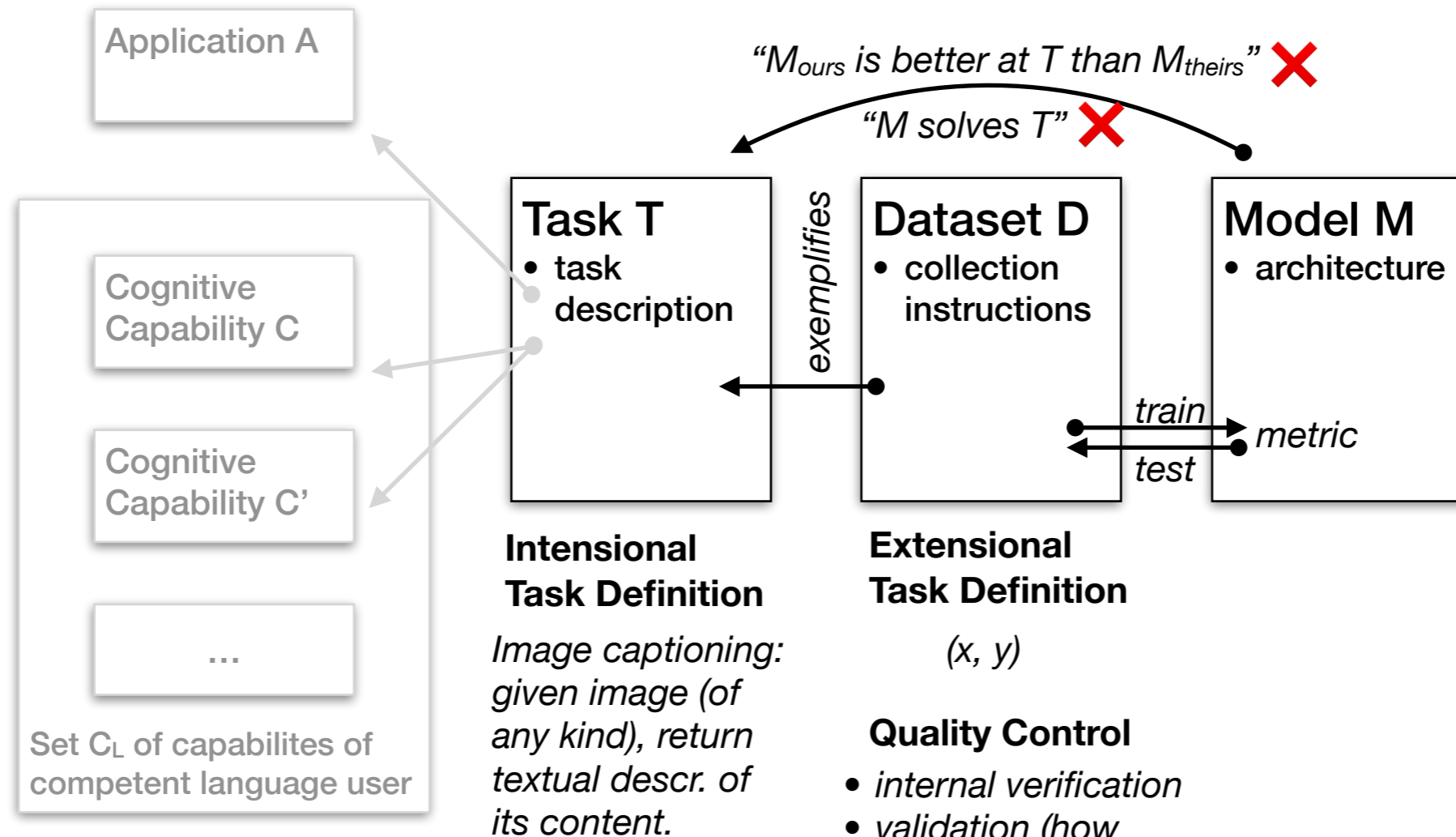
**Task Justification**

**Extensional Task Definition**

$(x, y)$

**Quality Control**

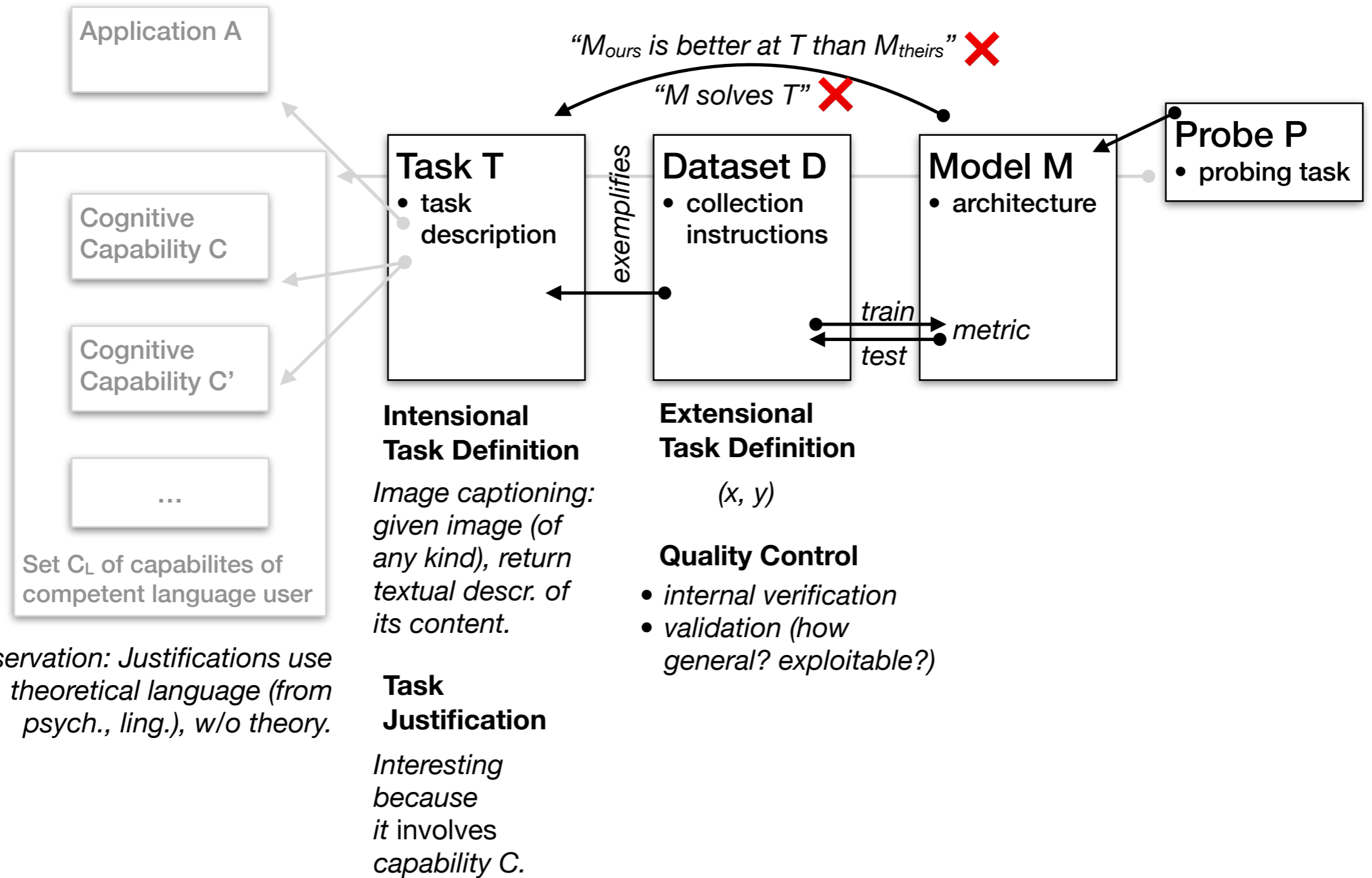
- internal verification
- validation (how general? exploitable?)



Observation: Justifications use theoretical language (from psych., ling.), w/o theory.

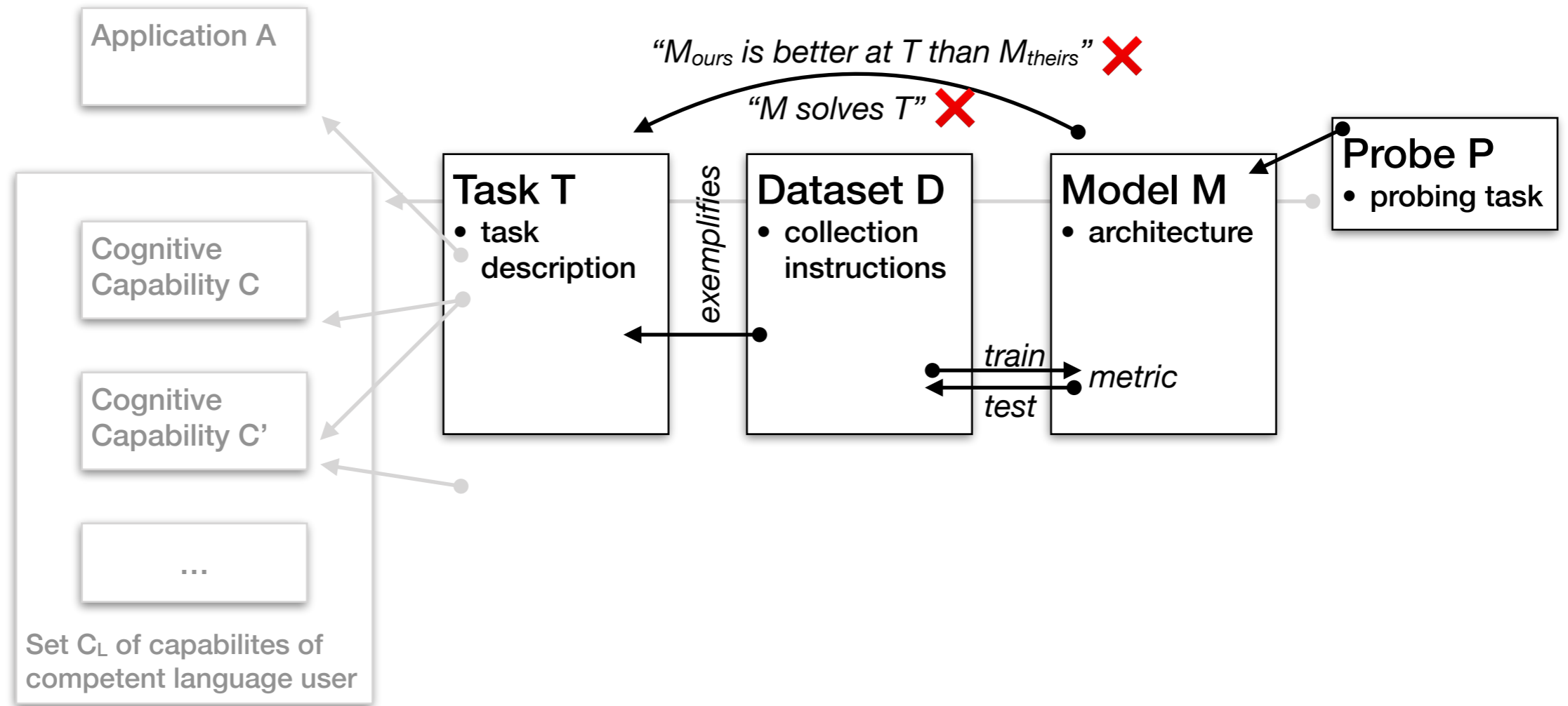
**Task Justification**

*Interesting because it involves capability C.*

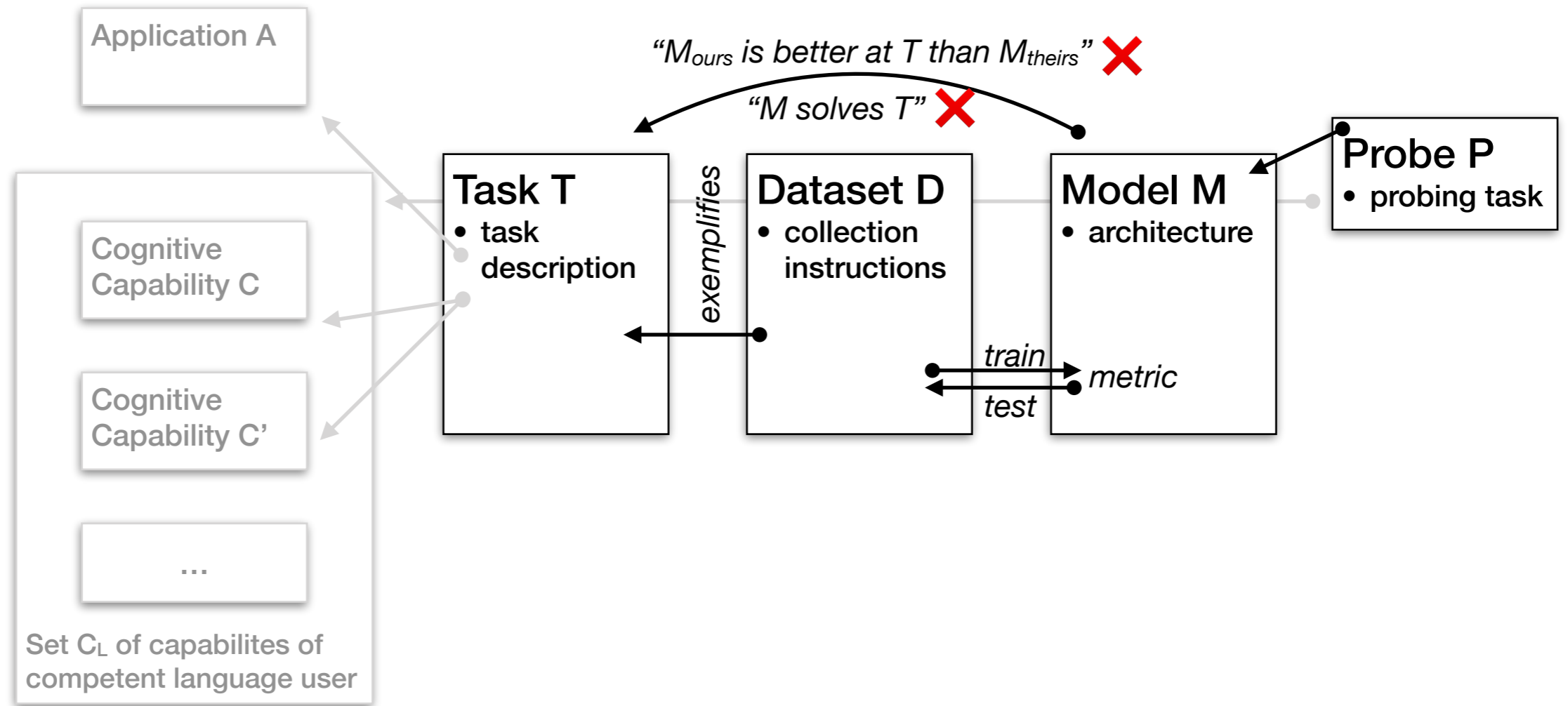


Observation: Justifications use theoretical language (from psych., ling.), w/o theory.



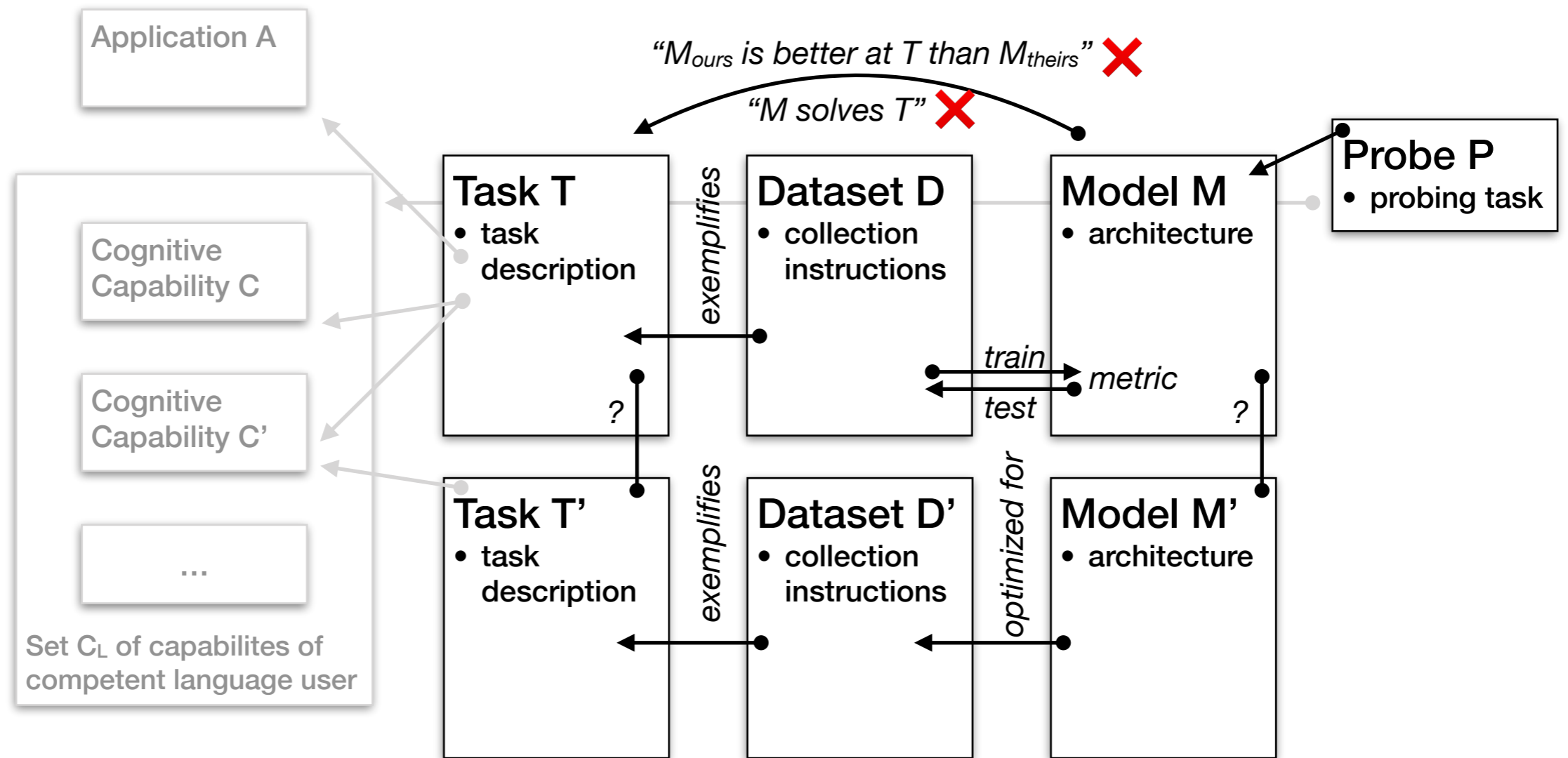


The Task / Dataset pair drives our research, perhaps more so than the Models. It should get appropriate attention.



The Task / Dataset pair drives our research, perhaps more so than the Models. It should get appropriate attention.

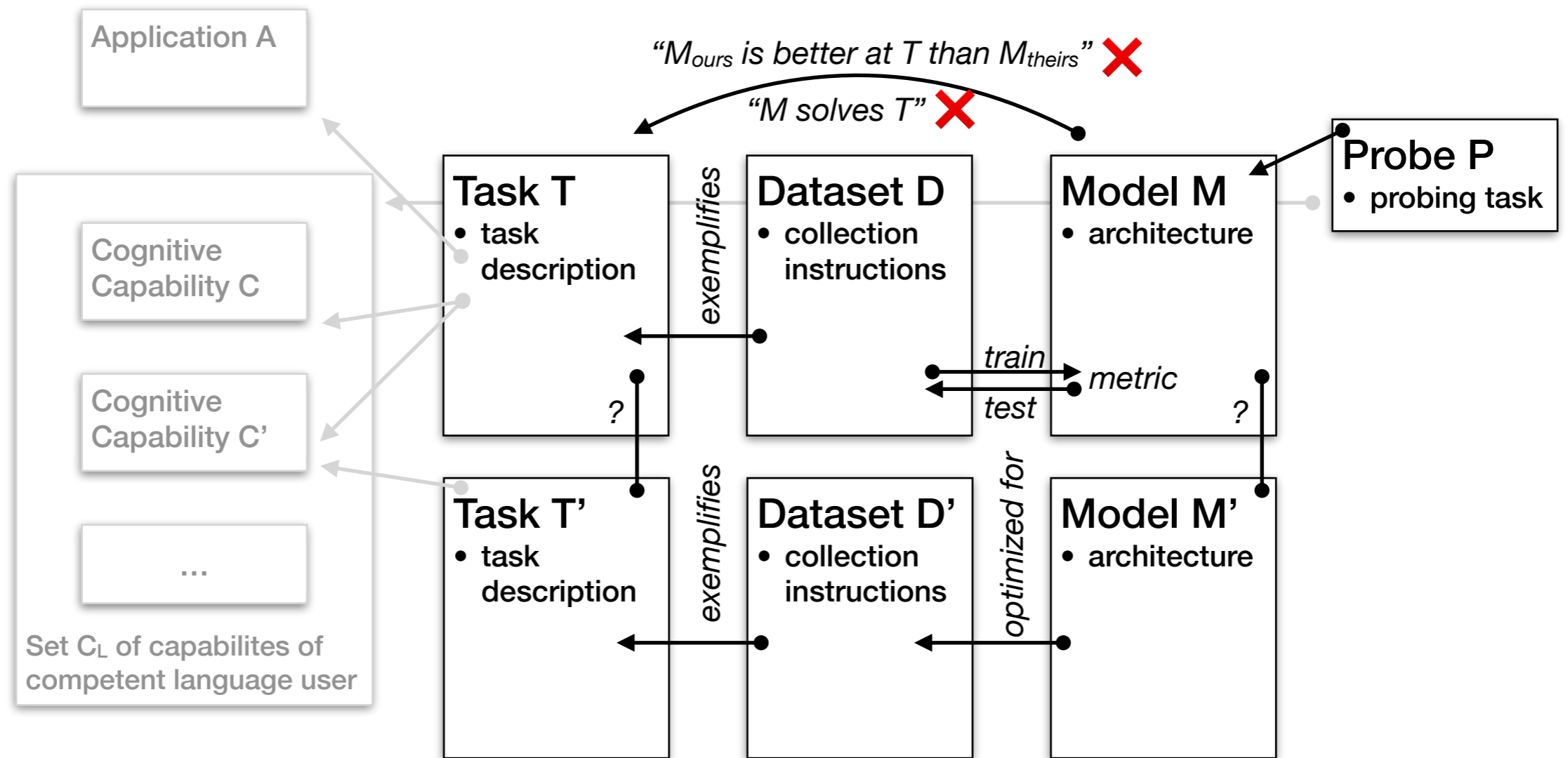
Can we find better ways to ground tasks in considerations about cognitive capabilities and the composition of the language faculty?  
(Or a principled argument for why we needn't care?)



The Task / Dataset pair drives our research, perhaps more so than the Models. It should get appropriate attention.

Can we find better ways to ground tasks in considerations about cognitive capabilities and the composition of the language faculty?  
(Or a principled argument for why we needn't care?)

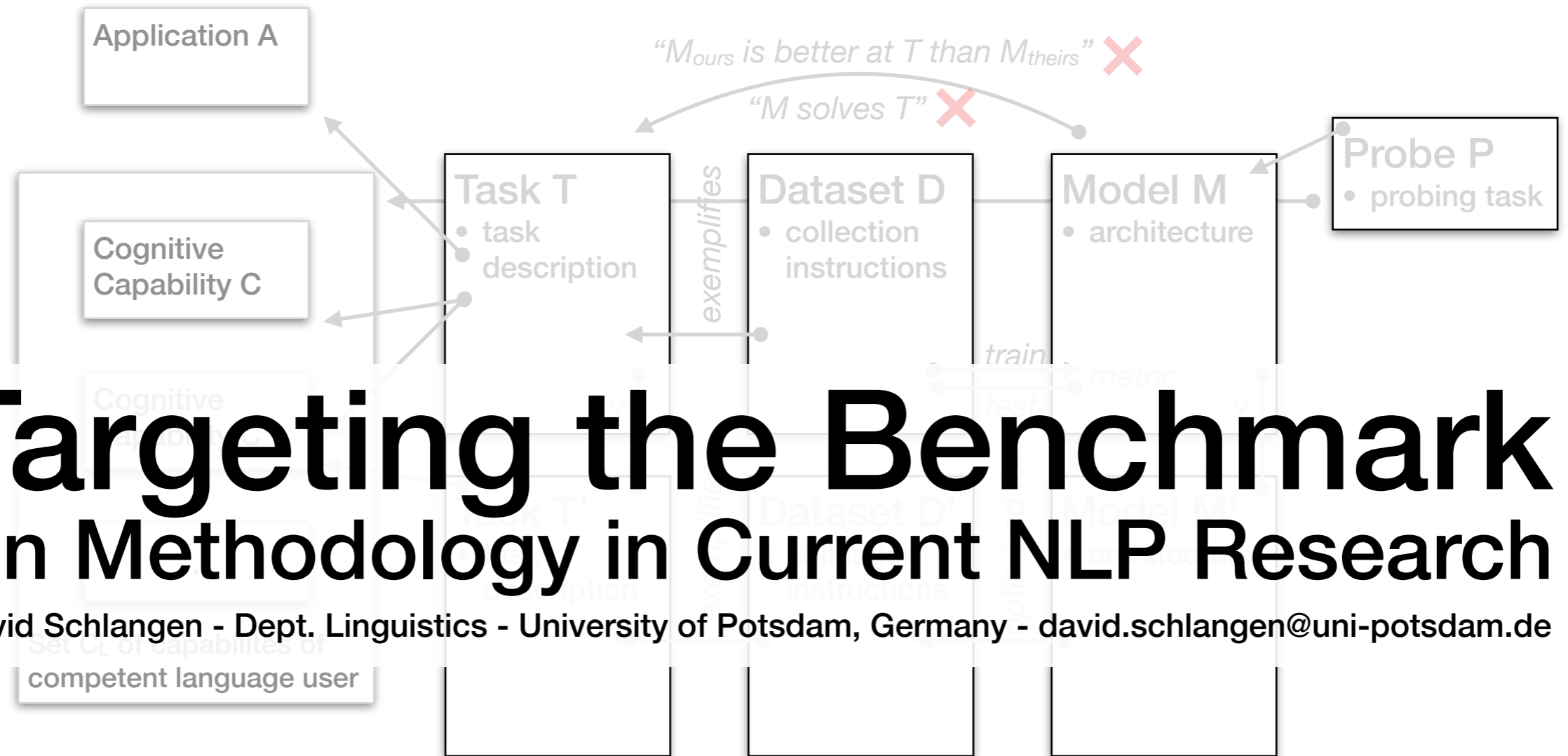
Can we make tasks *predictive*, so that performance of M at T tells us something about performance at T'?



The Task / Dataset pair drives our research, perhaps more so than the Models. It should get appropriate attention.

Can we find better ways to ground tasks in considerations about cognitive capabilities and the composition of the language faculty?  
(Or a principled argument for why we needn't care?)

Can we make tasks *predictive*, so that performance of M at T tells us something about performance at T'?



# Targeting the Benchmark On Methodology in Current NLP Research

David Schlangen - Dept. Linguistics - University of Potsdam, Germany - david.schlangen@uni-potsdam.de

The Task / Dataset pair drives our research, perhaps more so than the Models. It should get appropriate attention.

Can we find better ways to ground tasks in considerations about cognitive capabilities and the composition of the language faculty? (Or a principled argument for why we needn't care?)

Can we make tasks *predictive*, so that performance of M at T tells us something about performance at T'?