



Incentivising Pragmatism

Advanced methods for multi-agent communication

Mortimer von Chappuis
11 October 2023

Setup

Environment

- Self play
- partial observability
- fully cooperative

Symmetries

- Valuefunction
- Policy
- Belief update

Three levels of incentives

descriptive

$$\operatorname{argmax}_m I(m^{\text{self}}, b^{\text{self}})$$

epistemic

REDACTED

pragmatic

$$\operatorname{argmax}_m Q(b^{\text{other}}, u^{\text{other}} | m^{\text{self}}) - Q(b^{\text{other}}, u^{\text{other}})$$

reward shaping

$$r_t + \alpha \hat{r}_t - \beta |m_t|$$

Hypothesis

distractors

targets

novelty

familiar

descriptive



epistemic



pragmatic



REDACTED

Mutual Information

$$I(b, m) = \sum_{b \in B, m \in M} p(b, m) [\log(p(m|b)) - \log(p(m))]$$



$\pi(m|b)$

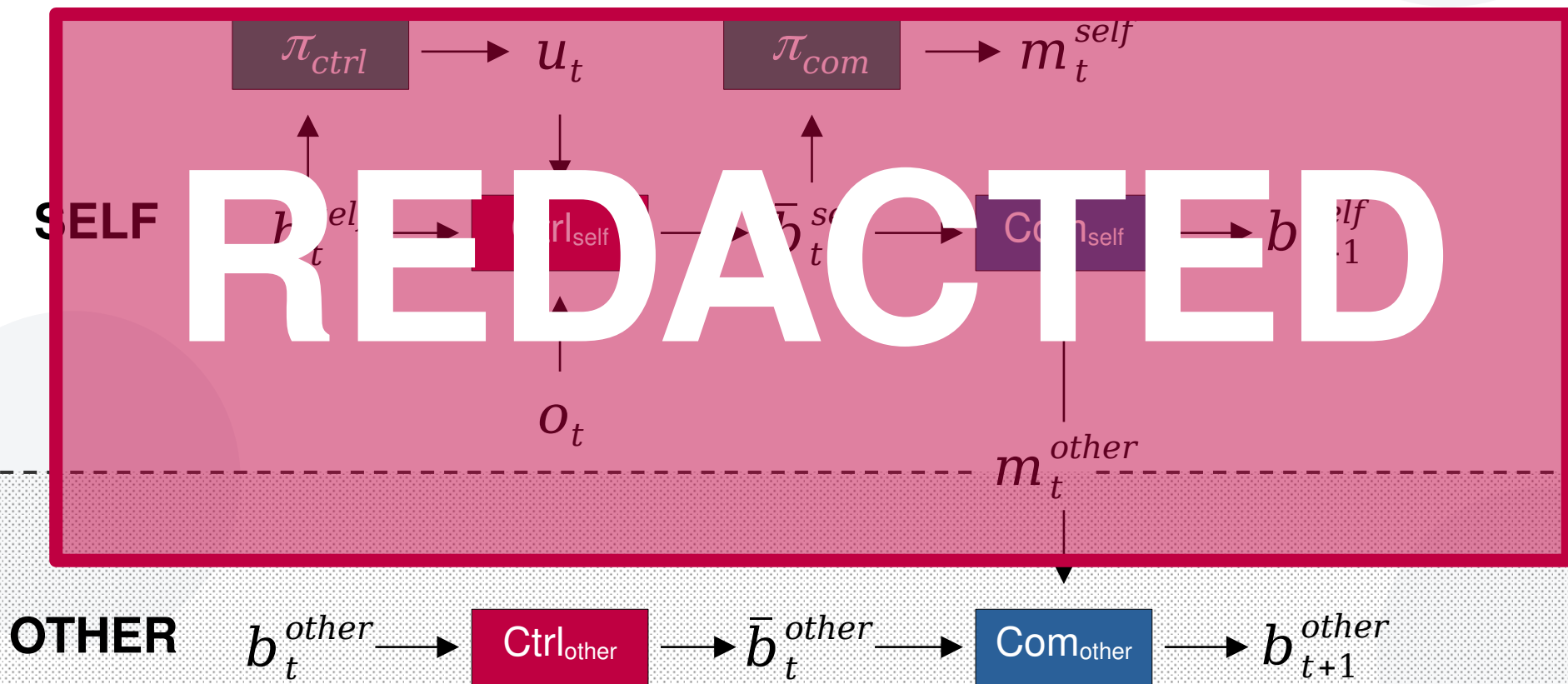
policy



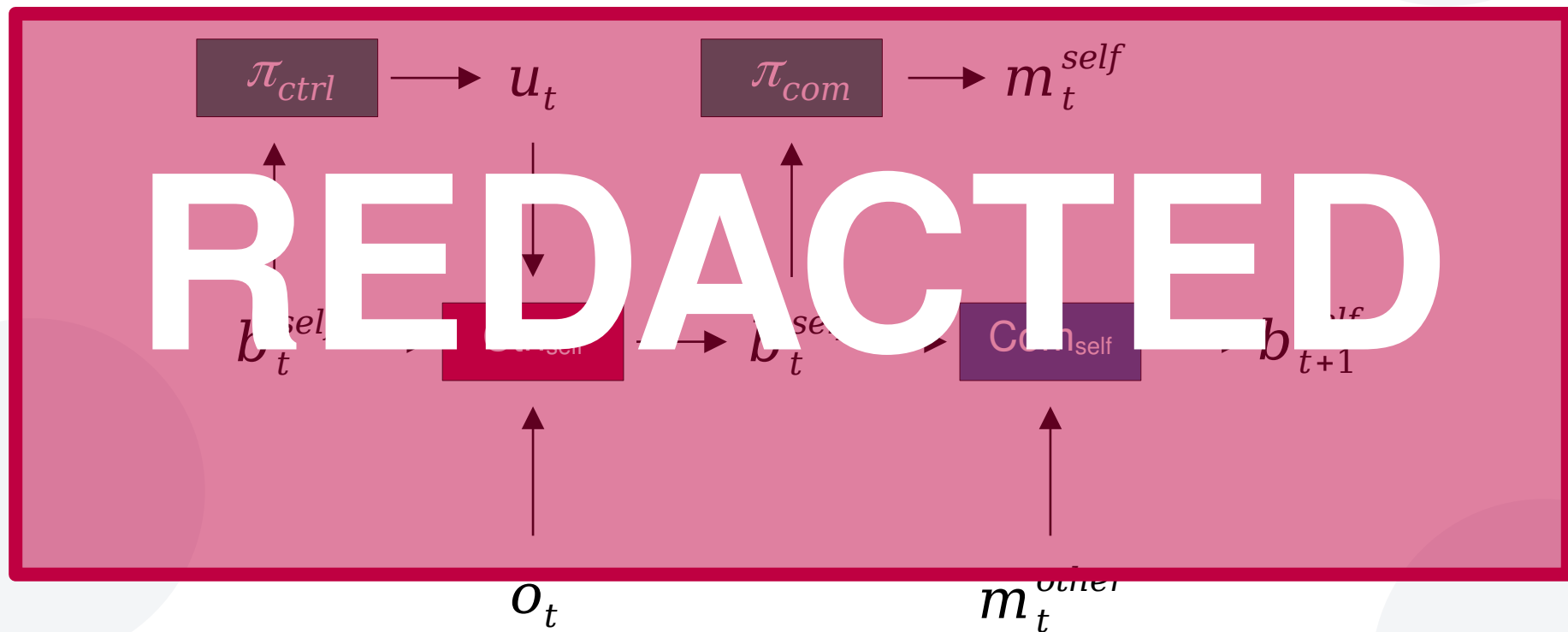
$f(m)$

language model

Architecture



Architecture — Self Belief

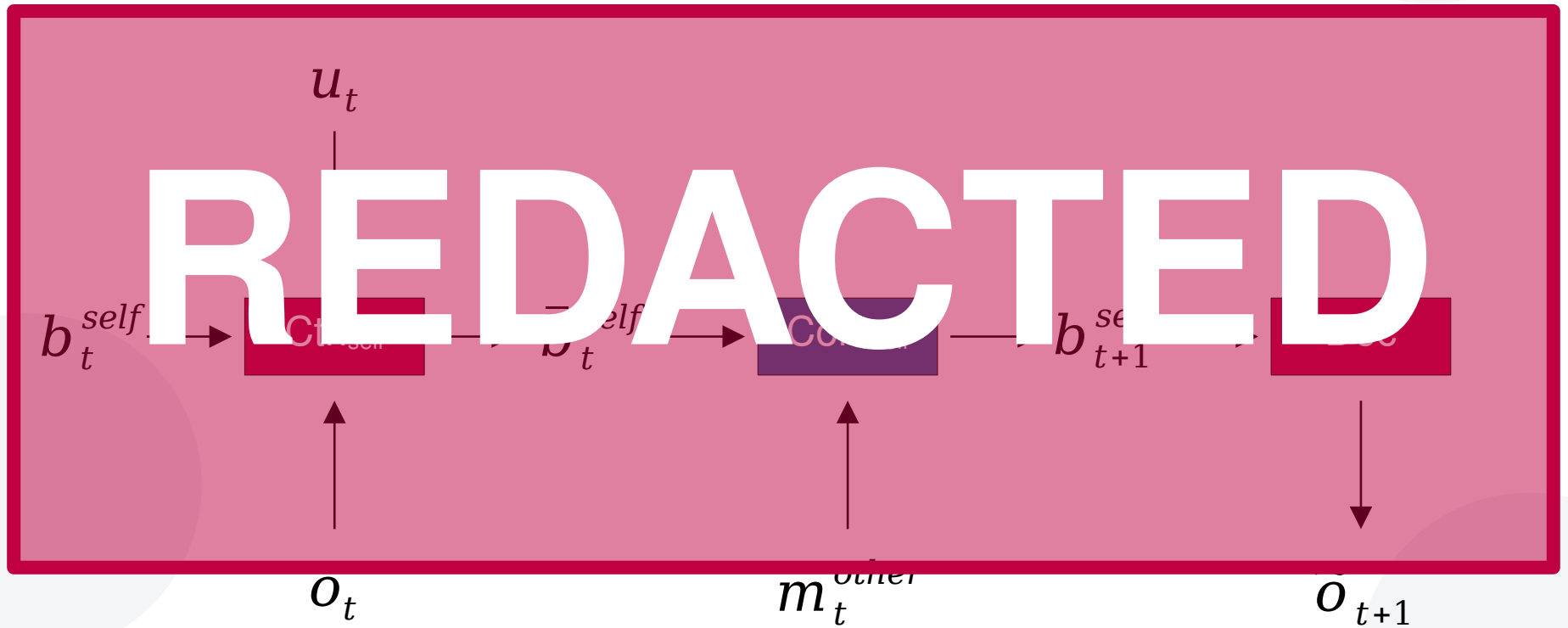


Architecture — Other Belief

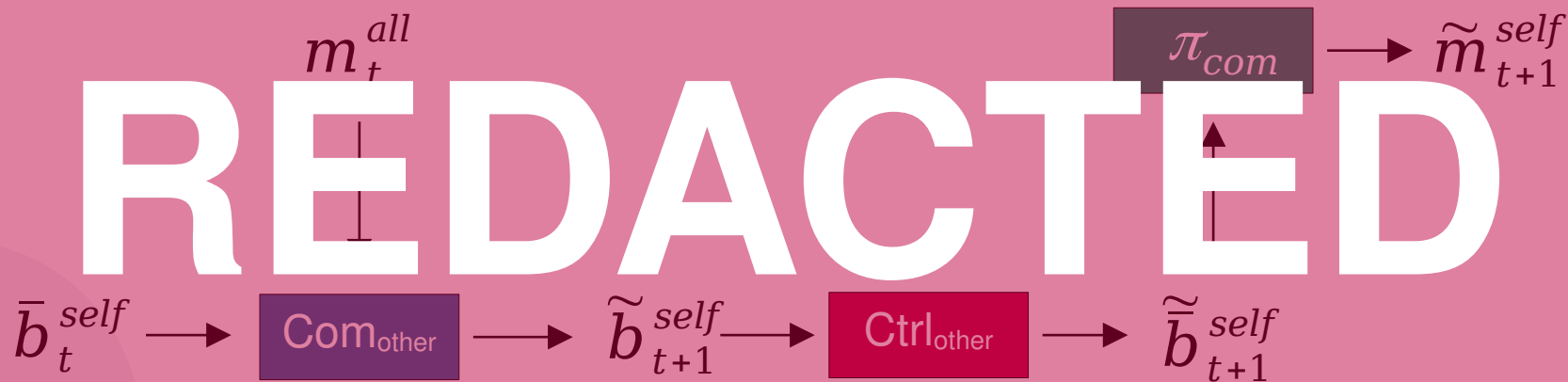
REDACTED



Training — Self



Training — Other



Summary

Networks

- actor policy network, policy (ctrl)
- critic Q function or advantage
- language model
- ctrl (self, other)
- com (self, other)
- dec

Training

- RL (actor, critic)
- language model
- belief self ($\text{Ctrl}_{\text{self}}$, Com_{self} , Dec)
- belief other ($\text{Ctrl}_{\text{other}}$, $\text{Com}_{\text{other}}$)

Thanks for your attention! — Questions?

Influences

- Cheap Talk Discovery and Utilization in Multi-Agent Reinforcement Learning
- Learning Attentional Communication for Multi-Agent Cooperation
- Mastering Diverse Domains through World Models