

AGI/ASI

- ANI (narrow AI): Good at one thing i.e.
(captioning img, playing chess), what we have now.

- AGI (Artificial general intelligence): can understand
learn and do most cognitive tasks that humans
can, across domains, with strong transfer + Autonomy.

- ASI (Artificial super intelligence): Far beyond
human performance across virtually all domains like
STEM.

↑
increasing
capability and
Generality
(i.e. tasks + domains)

- How to get AGI

- Scale + Data + Compute • Better architectures • Reasoning
- Agents • CoT • RAG • Multimodal • Distillation + reflection

- Why AGI is hard

- Alignment and Safety • Deceiving • Feedback / planning

o ASE = Beyond the best human experts at Everything but not AGI First

o AGI and ASE can do things independently with little to no
human feedback

NOTE: AGI ≠ consciousness → not req

Bigger ≠ AGI → Scaling helps but reasoning memory etc matters too

Boxing ≠ Enough → isolation helps but is not enough

Alignment by making it nice ≠ work! too probabilistic
cannot control output just like that

Safety: RLHF, REINFORCE