# SUPERINTELLIGENCE

HUMANS HAVE LIMITATIONS THAT AGI NOT HAVE
- SPEED DIFFERENCE BETWEEN NEURONS-TRANSISTORS
- BRAIN SIZE
- FASTER THAN EVOLUTION

MOST LIKELY WAY TO ACHIEVE IS THROUGH A COLLECTIVE AGI COMPOSED OF MULTIPLE AGIS

## TWO APPROACHES FOR AI
- TASK BASED APPROACH: DESIGN SPECIFIC WAY TO APPLY AI TO EACH TASK
- GENERALIZED-BASED APPROACH: UNDERSTAND NEW TASKS WITH LITTLE OR NO TASK-SPECIFIC TRAINING BY GENERALIZING FROM PREVIOUS EXPERIENCE (BOTTOM-UP)

REACHING SUPERINTELLIGENCE

EXCEED HUMAN PERFORMANCE IN TERMS OF PROCESSING SPEED AND SIZE OF NEURAL NETWORK

+

POTENTIAL OF REPLICATION, CULTURAL LEARNING AND RECURSIVE IMPROVEMENT

← WILL BECOME GREATER AS AIS BECOME MORE INTELLIGENT

"INTELLECT THAT GREATLY EXCEEDS THE COGNITIVE PERFORMANCE IN ALL DOMAINS OF INTEREST, BETTER THAN ALL OF HUMANITY COORDINATING GLOBALLY" ← BOSTROM

## 3 FACTOR THAT WILL BECOME MORE IMPORTANT AS AIS BECOME MORE INTELLIGENT
- REPLICATION:
  - AIS LESS CONSTRAINED THAN HUMANS. EASY TO CREATE A DUPLICATE OF IT
  - AGIS DECOMPOSE DIFFICULT TASKS IN SUBTASKS.
  - COLLECTIVE SUPERINTELLIGENCE: LARGE GROUP OF AGIS TO CARRY OUT MORE COMPLEX TASKS
- CULTURAL LEARNING:
  - SKILL TO ACQUIRE/SHARE KNOWLEDGE. COLLECTIVELY AGI SOLVE HARDER PROBLEM THAN INDIVIDUAL AGIS
- RECURSIVE IMPROVEMENT:
  - IMPROVE TRAINING PROCESSES TO DEVELOP THEIR SUCCESSOR

## POTENTIAL OF THE GENERALIZATION-BASED APPROACH IN HOW HUMANS DEVELOPED
- SKILL OF ABSTRACTION: EXTRACT COMMON STRUCTURE FROM DIFFERENT SITUATIONS = MORE EFFICIENT UNDERSTANDING
- COMMUNICATION SKILLS-THEORIES: SHARE OUR IDEAS

AGENCY: ABILITY OF AN AGENT TO HAVE ITS OWN GOALS:
- DESIGNED OBJECTIVES: GOALS THAT AN AI HAS BEEN DESIGNED TO ACHIEVE
- OWN GOALS: GOALS THAT AN AI WANTS TO ACHIEVE

CURRENT AI SYSTEMS: ACHIEVE DESIGN OBJECTIVES WITHOUT TRULY UNDERSTANDING WHAT ARE/ACTIONS TO ACHIEVE THEM

BOUNDED RATIONALITY: SYSTEM CAN TRY TO ACHIEVE A ROLE W/OUT TAKING THE BEST ACTIONS.

# GOAL AND AGENCY

3 WAY HOW AN AI COULD GAIN POWER
- AIS PURSUE POWER AS AN INSTRUMENTAL GOAL TO ACHIEVE OTHER GOALS
- AIS PURSUE POWER FOR ITS OWN SAKE
- AIS GAIN POWER WITHOUT AIMING FOR IT

THERE ARE INSTRUMENTAL GOAL THAT INCREASE THE CHANCES OF AN AGENT'S FINAL GOALS BEING REACHED
- SELF-PRESERVATION
- RESOURCE ACQUISITION
- SELF-IMPROVEMENT

## GOAL-DIRECTED
1) SELF-AWARENESS: UNDERSTAND THAT AN AGENT IS PART OF THE WORLD AND THAT ACTIONS HAVE CONSEQUENCES
2) PLANNING: MAKE DECISIONS WRT THE OUTCOME VALUE
3) CONSEQUENTIALISM: DECIDE BEST PLANS WRT CONSEQUENCES OF ACTION
4) SCALE: TAKE INTO ACCOUNT EFFECTS OF PLANS
5) COHERENCE: REMAIN INTERNALLY UNIFIED IN IMPLEMENTING THE SINGLE BEST ACTION
6) FLEXIBILITY: ADAPT PLANS FLEXIBLY

## WHAT IS MISSING
1) TRAINED ON 3RD-PERSON DATA, NO 1ST-PERSON PERSPECTIVE
2) LIMITED TYPES OF PLAN CONSIDERED. NOT ALWAYS EXTENSIVE PLANS
3) SUBJECTIVE PREFERENCES ABOUT ACTIONS WRT CONSEQUENCES
4) TRAIN IN SMALL SCALE ENVIORMENT/NO GENERALIZATION ABILITY
5) INTERNAL CONFLICT (LIKE HUMANS)
6) 1 INITIAL PLAN, NO ADAPTATION TO NEW SITUATIONS/RETHINK PLANS