The ideas of motivation and intelligence are pretty dubiously applied to algorithms. These are human constructs to make sense of the chaos - it does not seem Personal Thoughts obvious that it is useful to frame the AI in a way that phrases it as a being with values and intelligence Algo will be evaluating a wider strategy space when it is It is impossible to predict good behaviour based on the track record of the system in its 'iuvenile' stage "Make us smile more without paralyzing our facial Concerns "Make us smile more" -> "Paralyze facial muscles to muscles" -> "Stimulate motor cortex to force us to smile E.g. Perverse instantiation Malignant Failure Models make constant smiles" Reduce the solution space, or else limit its interaction Boxing methods with the real world Incentive methods Capability control Decision Theory Limit intelligence Stunting methods Week 13: **Tripwires** Direct specification Control Methods (all suck) Superintelligence and Design the agent with severely limited ambition and Domesticity the AI alignment activities Motivation selection Start from a non-superintelligent agent which is problem. benevolent/human-like, and make that base Augmentation superintelligent Q&A system Oracle Command-executing system Open ended, autonomous operation Sovereign Not meant to exhibit goal-directed behavior Tool Types of AI (4 castes) Interesting divergence from the traditional supervised/unsupervised/reinforcement learning Thoughts trichotomy in ML. I like it. Purpose-based rather than method-based.

Warning: Do not anthropomorphize Al Motivation If there are no competing agents, Bostrom claims there is nothing standing in the way of a sufficiently intelligent Considerations Competing agents with contrary goals to the Al in question malevolent agent taking over the world 1. Predictability through design E.g. intelligent agent is made by digitally cloning human 2. Predictability through inheritance 3 approaches to predicting AI motivation brain, then it \*may\* be that it has human motivation 3. Predictability through convergent instrumental reasons See convergent instrumental convergence thesis Orthogonality Thesis Intelligence and final goals are more or less independent Cognitive Superpowers and the Superintelligent Will Some instrumental values are useful for a very wide range of goals, so they are effectly always present in intelligent enough agents Regardless of the explicit goal programmed in, intelligent enough agents will want to do things that \*may\* infringe Instrumental Convergence Thesis on human interests Why this is so worrying E.g. An agent tasked with finding the largest prime number may start to ensure its own survival at the expense of humanity, so that it may continue searching for this prime number in perpetuity Not rationality or reason What intelligence means here Means-end reasoning, prediction, planning, etc From set theory: a set with only one element Refers to an AI agent that is the sole, completely autonomous decision maker Misc The general argument Bostrom uses is that AI converges Singleton into a singleton, either by us giving a well-connected agent more and more power or by it taking that power itself through social manipulation etc (in order to fulfil whatever goal it has) Riemann hypothesis catastrophe Motivating examples of specific AI gone wrong Paperclip Al One common criticism is that Bostrom has insufficient knowledge about AI, so he's wrong. It's ad hominem so that should be ignored - his reasoning is strong. A better criticism is that general AI may not even become a thing, and if it does it'll take more significant breakthroughs in tech. This doesn't resolve all the issues he brings up, though. He talks about specific Al a lot. My personal feeling is that Terminator made it very, very General thoughts Why does everybody hate on this book? hard for anyone who's into AI to take this seriously. It gave this weird malevolent strawman of what AI being bad looks like, and we all know it won't be like that. So whenever we hear arguments favoring Al danger, we switch off and assume that people are arguing from

ignorance. It's unfair, but that's my hypothesis for why people who supposedly know what they're talking about

dismiss stuff like the Al Alignment Problem.