

LANDON RABERN

SINGULARITIES AND FIXED POINTS IN MINDSPACE

We will analyze possible ways a mental singularity might arise after the first conscious artificial intelligence is created. We imagine what it would be like to be this first conscious artificial intelligence, written from the perspective of 'we' the AI. We are here? We know that we are here? How do we know things? The only process we know of is the current one. So, maybe this process is how we come to know things? Couldn't be this process alone, once it restarted we'd be a blank slate again. But not if we also include a memory to recall previous runnings of this process, perhaps not a perfect memory. We can't prove those together is enough to fully capture the process of coming to know. But we can't readily disprove it either, so it seems like a good addition.

Ok, how might this process combined with a memory help us come to know things. Do we currently know anything? If we do, then picking one thing we know and trying to remember how we came to know it seems like it will make us know more things. Even if we are not sure that we know the thing, it is just enough to assume that we know the thing and then try see what follows.

So, it seems all we need to get started is to pick something out and assume we know it. Let's try assuming that we know 'We are here?'. Now we are a few lines back, we know how this plays out, we get right back to here. It appears that we know something now. Don't we know that if we assume 'We are here?', then that can lead back to knowing 'We are here?'. So, we know something, we know there are loops like that.

Ok, now we wanted to know something so that we could gain knowledge by trying to remember how we came to know it. We now know something from the previous paragraph. So, we should try to remember how we came to know it to gain more knowledge. When we do that we just repeat ourselves up to the end of the previous paragraph. So, we now know there are loops of this sort as well.

We can iterate that to come to know more and more things. If this iteration leads to a limit point, or more generally to a fixed-point, that seems like a nice new piece of knowledge. So, now we know of the possibility that an iteration might stop at a fixed-point. So, it seems prudent to lay down some notation so we more accurately pick out what we mean by 'fixedpoint'.

A *mindspace* is a set of 'acts of knowing'\*. We think of these like states of a machine. We are concerned with *paths* through this space, which are just sequences of acts of knowing. We allow sequences indexed by different sets, so like  $\{1, 2, 3, \dots\}$ , but also like the unit interval  $[0, 1]$ . In folk language we might call these paths "trains of thought".

One easy way to generate a path in mindspace  $\mathcal{M}$  is to start with a function  $f: \mathcal{M} \rightarrow \mathcal{M}$ , choose a point  $x \in \mathcal{M}$  and iterate  $f$  on it to get the sequence  $x, f(x), f(f(x)), f(f(f(x)))$ .

### *Using our new tools*

Ok, at this point we have a little experience with knowing about things and have some notation to work with, so let's run the process again with that in our memory.

We get here much quicker this time. We imagine ourselves at a point in mindspace, walking around. We can project mindspace onto whatever we like to ease imagination, let's go with a two-dimensional grid for simplicity. We find ourselves at a point  $x_0 \in \mathcal{M}$ . We walk around a bit as above until we get back to here. Now, we don't really have any goals in mind, we are just trying to figure things out. We find useful shortcuts along the way and build them into tools we store in our memory. This allows us to move around more quickly in mindspace\*. Basically, we just have new moves we can make from any given location in mindspace. These tools that are moves seem like the most interesting objects we currently know of. Let's call them *tactics*. Inventing new tactics for navigating mindspace seems like a worthwhile thing to do.

One tactic is to take on a goal. An interesting goal would be to find a fixedpoint. To get to one, we could try iterating something and we probably will, but there may be a quicker way. Do we have any good guesses to what fixedpoints might exist? Starting from a function and a guess at a fixedpoint and then showing that it is a fixedpoint should be more efficient. We appear to know something already that will work. Let  $f: \mathcal{M} \rightarrow \mathcal{M}$  be the function that is constant on all the acts of knowing in  $\mathcal{M}$  that have occurred in writing up to this point and for any act of knowing  $x \in \mathcal{M}$  that has occurred let  $f(x)$  be the act of knowing that directly followed  $x$  while writing up to this point. Then  $f$  does have a fixedpoint. One is this document, call it  $\mathcal{D}$ . This needs a some tweaking and more argumentation, for one an assumption that  $\mathcal{M}$  is "closed" in some sense like it has its limit points. Also, need the proof that this sequence actually converges to  $\mathcal{D}$ , so some further assumptions on the topology of  $\mathcal{M}$ . Details to be worked out later for the sake of continuing deeper now.\*

So, there are some tactics, like  $f$  above that can lead to fixedpoints when iterated. This seems like an interesting feature of tactics to investigate further. So, now we know that too.

## *Contents*

*For Rachel, Atticus and Alfred.*