On attempting to reify a few of the things we may mean by "consciousness" with code

Josh Joseph, Dhaval Adjodah, Joichi Ito Massachusetts Institute of Technology



Why attempt to reify philosophy with code

- Lots of the words philosophers use describing aspects of consciousness tends shows up in CS/AI research
 - Mind, awareness, imagination, reasoning, consciousness, etc.
- Our intuition is CS/AI could benefit from a deeper understanding of philosophy
 - But telling people to read more books/papers is not how to make this happen
 - So let's try to do it with code!
- Possibly benefit philosophy by bringing code-style concreteness
 - (TBD, will let the philosophers in the room speak to this!)

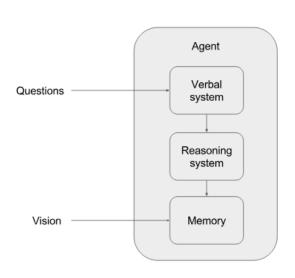
• (Disclaimer: our backgrounds are CS/AI)

Reifying philosophy with code

 Muehlhauser, Shlegeris: A Software Agent Illustrating Some Features of an Illusionist Account of Consciousness

An agent that observes the world and uses a theorem prover to answer

questions asked of it



```
Q: What's 2 + 2?

4

Q: Suppose there are two agents Bob and Jane, do they have the same qualia associated with every color? Both that statement and its negation are possible.

Q: For all y, does there exist an x such that x = y + 1?

Yes.

Q: For all two agents, do they see colors the same?
Both that statement and its negation are possible.

Q: Are your memories at timestep 0 and 1 of the same color?

Yes.

Q: Are you seeing the same color now as you saw at timestep 0?

No.

Q: Is it possible for an agent to have an illusion of red?

Yes.

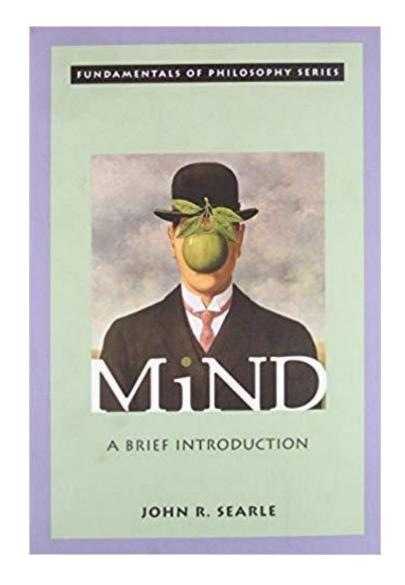
Q: Is it possible for you to have the illusion that Buck is experiencing a color?

Yes.

Q: Is it possible for Buck to have an illusion that he is having the experience of redness?

No, that's impossible.
```

Reifying philosophy with code



Reifying philosophy with code

- Searle's view of the relationship between consciousness and brain states
 - Consciousness is causally reducible to brain states but consciousness is ontologically irreducible to brain states
 - ...what does that mean?
- Generally is some confusion
 - Enough disagreement that Searle wrote the paper: "Why I'm Not a Property Dualist"
- Let's unpack this with code!

What we're not doing

- Not trying to
 - Propose a cognitive architecture
 - Propose a new AI or machine learning algorithm
 - Claim that the software agent is conscious
 - Convince anyone these are the correct/best/most useful definitions of consciousness or brain states
 - Convince anyone Searle is right or wrong

What we're trying to do

- Trying to create a software agent that is consistent with Searle's view on consciousness
 - (or at least a simplified version of Searle's view)
- (Hopefully) gain a bit deeper understanding of what we may mean by consciousness, brain states, causal reduction, and ontological reduction along the way

Software Engineering, 101

- Requirements what the system must do
- Design how will we build the system to meet the requirements
- Implementation building the system consistent with the design

- Consciousness is causally reducible to brain states
- Consciousness is ontologically irreducible to brain states

Brain state

- The full physical-chemical state of the brain and nervous system
- Third person, objective

Internal state

- Representations, goals, rewards, observations, actions, etc.
- Subjective

Mental state

- Beliefs, desires, thoughts, perceptions, emotions, knowledge, etc.
- First person, subjective

Conscious mental state

- A mental state in which it is "something it's like to be in"
- First person, subjective character of experience, phenomenal

Searle's view

- Consciousness is causally reducible to brain states
- Consciousness is ontologically irreducible to brain states

...simpler

- Conscious mental states are causally reducible to brain states
- Conscious mental states are ontologically irreducible to brain states

• ...simpler

- Mental states are causally reducible to brain states
- Mental states are ontologically irreducible to brain states

• ...simpler

- Internal states are causally reducible to brain states
- Internal states are ontologically irreducible to brain states

- Searle's view
 - Consciousness is causally reducible to brain states
 - Consciousness is ontologically irreducible to brain states
- V2
 - Conscious mental states are causally reducible to brain states
 - Conscious mental states are ontologically irreducible to brain states
- V1
 - Mental states are causally reducible to brain states
 - Mental states are ontologically irreducible to brain states
- V0
 - Internal states are causally reducible to brain states
 - Internal states are ontologically irreducible to brain states

- V0
 - Internal states are causally reducible to brain states
 - Internal states are ontologically irreducible to brain states

Phenomena of type A are ontologically reducible to phenomena of type B if and only if A's are nothing but B's

Class-instance distinction



C Wine

Images from:

Class-instance distinction



C Wine

https://protege.stanford.edu/publications/ontology Case of wine https://protege.stanford.edu/publications/ontology Case of wine https://www.researchgate.net/figure/Owl-Viz-view-of-course-ontology fig1 261339041

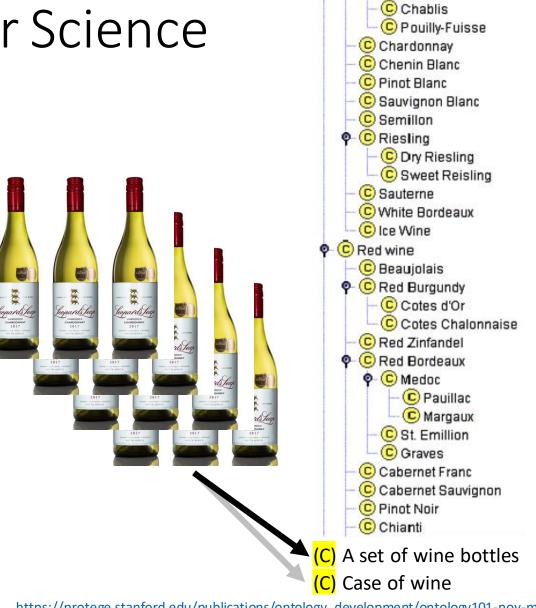
Class-instance distinction



C Wine C White wine C Rose wine C Red wine C White Burgundy C Chenin Blanc C Chardonnay C Pinot Blanc C Sauvignon Blanc C Ice Wine C White Zinfandel C Beaulolais C Red Burgundy C Red Zinfandel C Pauillac C Margaux C St. Emillion C Graves C Red Bordeaux © Sauterne C Cabernet Franc C Cabernet Sauvignon C Medoc © Semillon C Pinot Noir C Chianti C Petite Syrah C Sancerre C Muscadet C Port C Sweet Reisling C Chablis C Dry Riesling (C) A set of wine bottles

https://protege.stanford.edu/publications/ontology Case of wine https://protege.stanford.edu/publications/ontology Case of wine https://www.researchgate.net/figure/Owl-Viz-view-of-course-ontology fig1 261339041

- Class-instance distinction
- Type-token distinction
 - "They drive the same car"
 - They drive the same car type
 - (a Toyota)
 - They drive the same car token
 - (the 2003 Toyota Corolla with VIN: 2QFBORHE4KP911561)
- Representing tokens of one type as tokens of another type

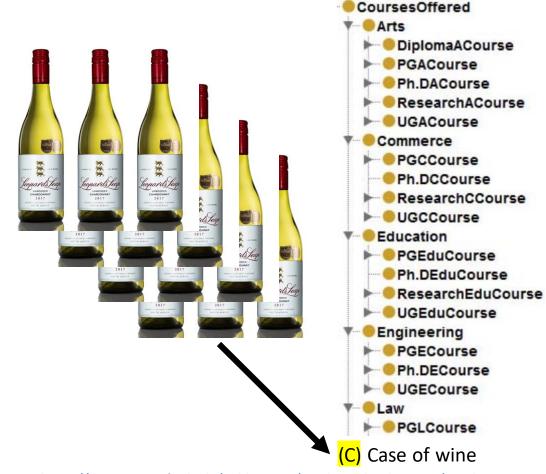


C) Wine

- C White wine

White Burgundy

- Class-instance distinction
- Type-token distinction
 - "They drive the same car"
 - They drive the same car type
 - (a Toyota)
 - They drive the same car token
 - (the 2003 Toyota Corolla with VIN: 2QFBORHE4KP911561)
- Representing tokens of one type as tokens of another type



Images from:

- V0
 - Internal states are causally reducible to brain states
 - Internal states are ontologically irreducible to brain states

Phenomena of type A are ontologically reducible to phenomena of type B if and only if A's are nothing but B's

- V0
 - Internal states are causally reducible to brain states
 - Internal states are ontologically irreducible to brain states

Phenomena of type A are causally reducible to phenomena of type B if and only if:

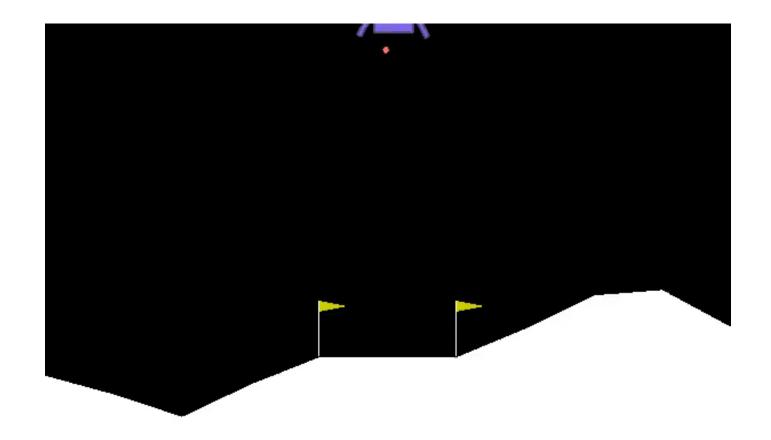
- the behavior of A's are entirely causally explained by the behavior of B's
- A's have no causal powers in addition to the powers of B's

Requirements, VO

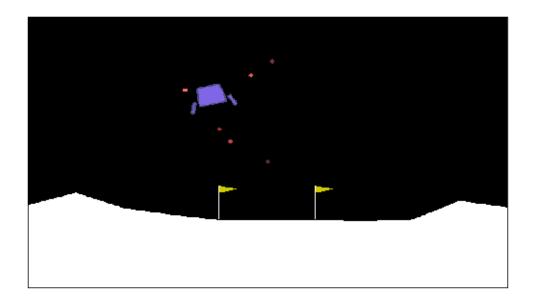
- Internal states are causally reducible to brain states
- Internal states are ontologically irreducible to brain states

Design decisions

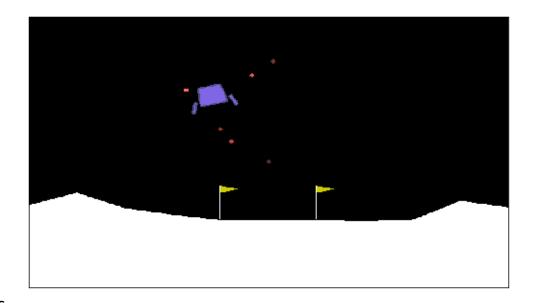
• OpenAI's LunarLander benchmark environment



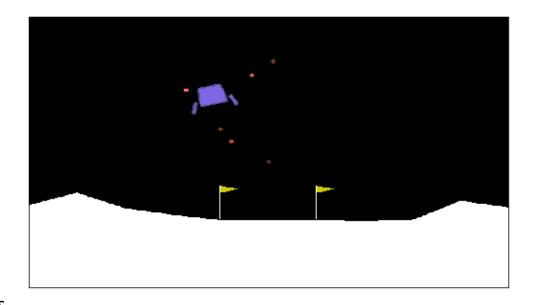
- Design decisions
 - Environment and the agent's "physical" form

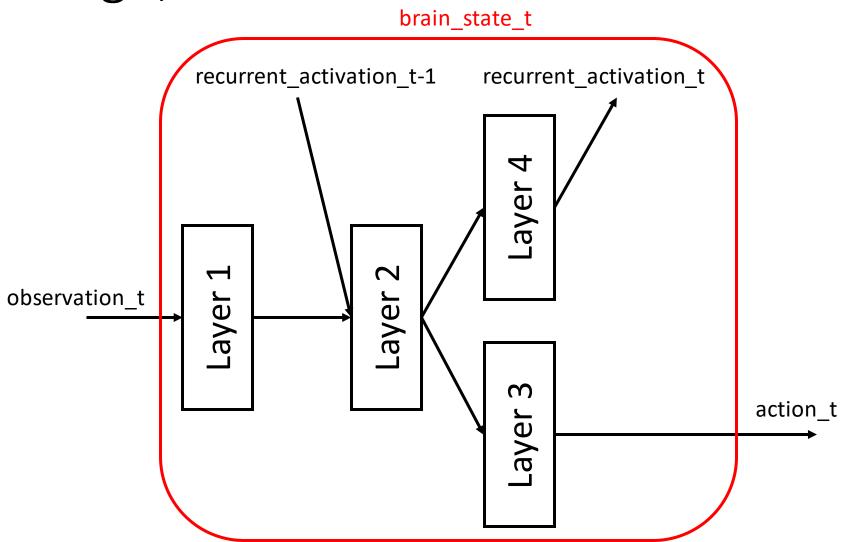


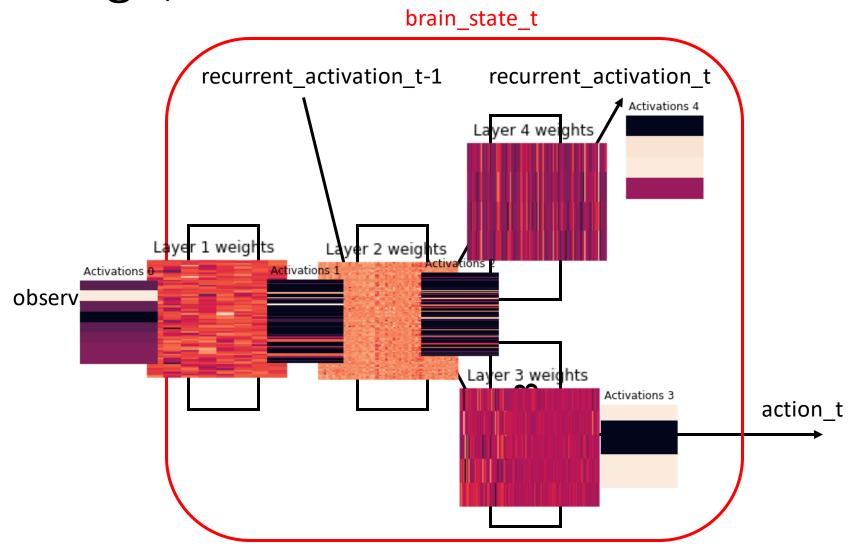
- Design decisions
 - Environment and the agent's "physical" form
 - Internal state of the agent
 - Beliefs about itself relative to semantically important regions
 - Left of the flags, right of the flags, high above the ground, close to the ground, falling too fast

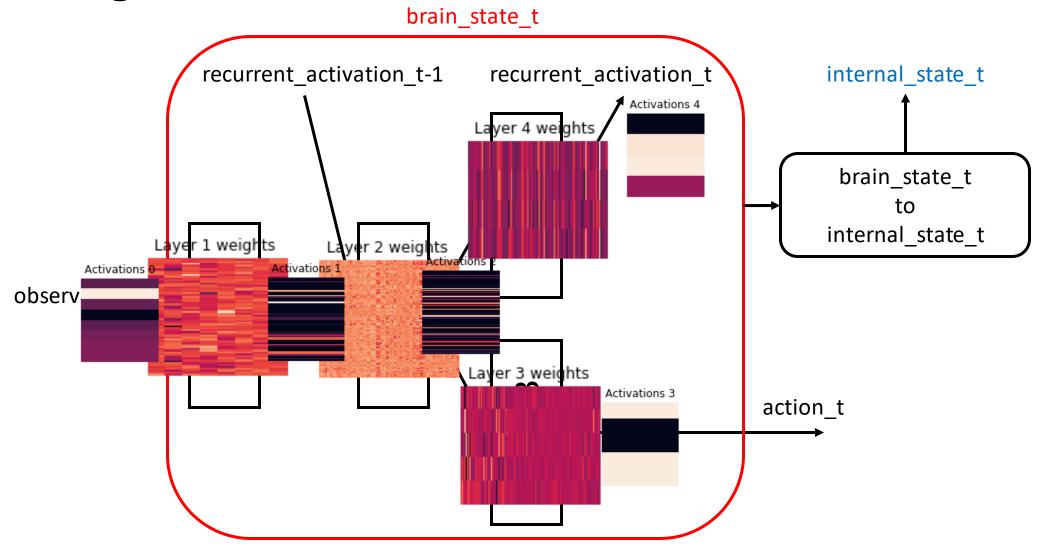


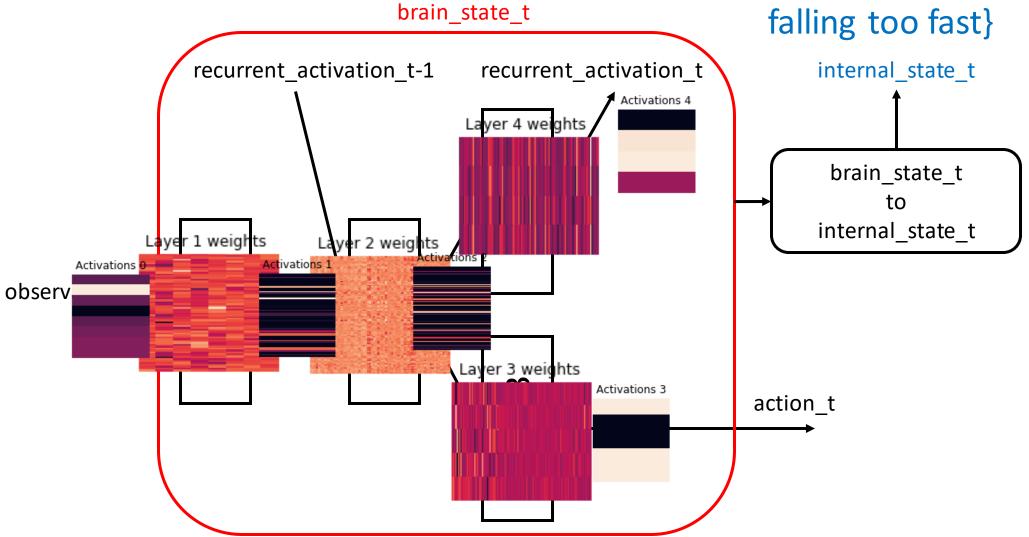
- Design decisions
 - Environment and the agent's "physical" form
 - Internal state of the agent
 - Beliefs about itself relative to semantically important regions
 - Left of the flags, right of the flags, high above the ground, close to the ground, falling too fast
 - Brain state of the agent





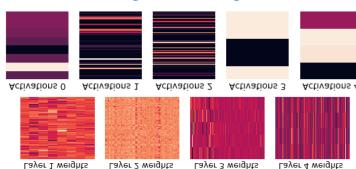


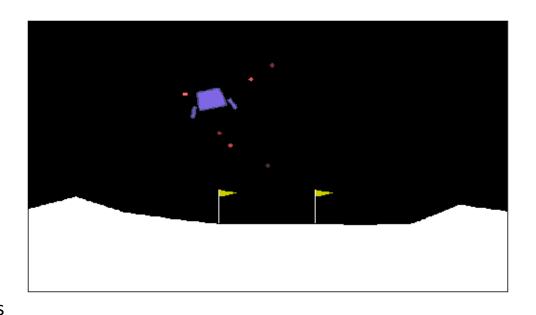




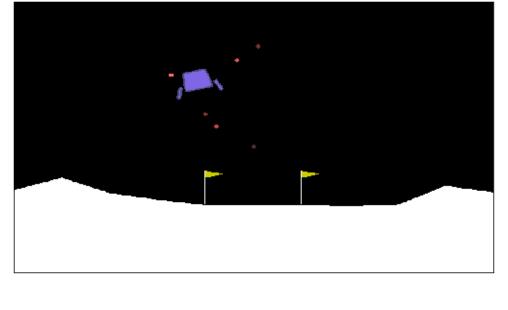
{High above the ground, right of the center falling too fast}

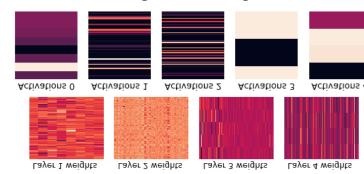
- Design decisions
 - Environment and the agent's "physical" form
 - Internal state of the agent
 - Beliefs about itself relative to semantically important regions
 - Left of the flags, right of the flags, high above the ground, close to the ground, falling too fast
 - Brain state of the agent



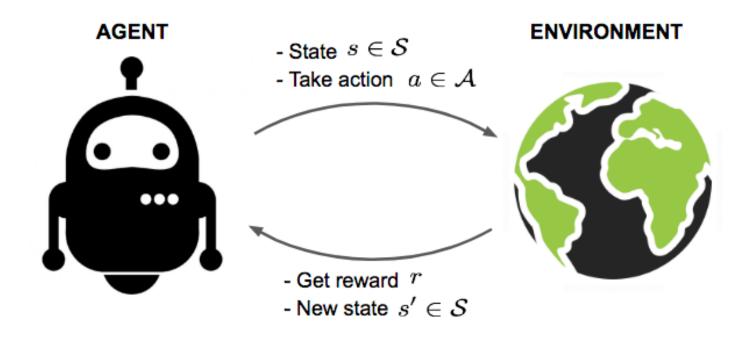


- Design decisions
 - Environment and the agent's "physical" form
 - Internal state of the agent
 - Beliefs about itself relative to semantically important regions
 - Left of the flags, right of the flags, high above the ground, close to the ground, falling too fast
 - Brain state of the agent
 - Our ontology
 - Layer weights of the neural network
 - Connectivity of the neural network
 - Activations of the neural network at time t
 - The agent's observation at time t
 - The agent's action at time t
 - The position and velocity of the agent at time t
 - A region the agent believes it's in
 - Brain state at time t (set of layer weights, activations, and connectivity)
 - Internal state at time t (set of regions the agent believes it's in)





Reinforcement learning



Implementation, VO

- Jupyter notebook time!
 - http://localhost:8888/notebooks/notebooks/TSC-2019.ipynb
 - https://github.com/Josh-Joseph/tsc-2019/blob/master/notebooks/TSC-2019.ipynb

- V0
 - Internal states are causally reducible to brain states
 - Internal states are ontologically irreducible to brain states

- V0
 - Internal states are causally reducible to brain states
 - Internal states are ontologically irreducible to brain states

Phenomena of type A are causally reducible to phenomena of type B if and only if:

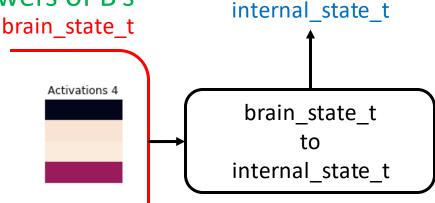
- the behavior of A's are entirely causally explained by the behavior of B's
- A's have no causal powers in addition to the powers of B's

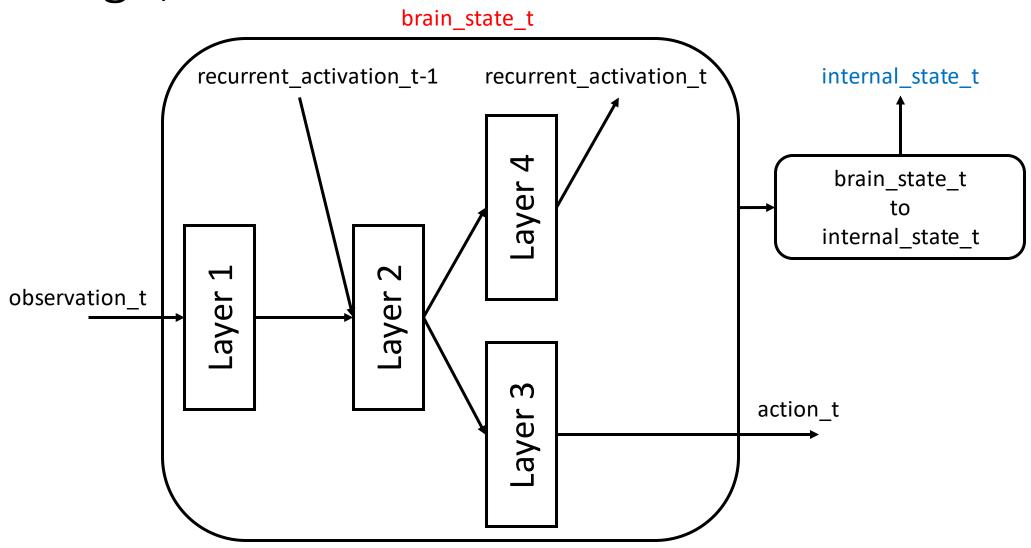
- V0
 - Internal states are causally reducible to brain states
 - Internal states are ontologically irreducible to brain states

Phenomena of type A are causally reducible to phenomena of type B if and only if:

- ✓ the behavior of A's are entirely causally explained by the behavior of B's
- A's have no causal powers in addition to the powers of B's

{High above the ground, right of the center falling too fast}



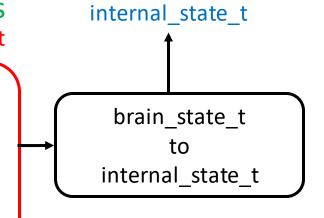


- V0
 - ✓ Internal states are causally reducible to brain states
 - Internal states are ontologically irreducible to brain states

Phenomena of type A are causally reducible to phenomena of type B if and only if:

- √ the behavior of A's are entirely causally explained by the behavior of B's
- ✓ A's have no causal powers in addition to the powers of B's brain state t

{High above the ground, right of the center falling too fast}



Activations 4

- V0
 - ✓ Internal states are causally reducible to brain states
 - Internal states are ontologically irreducible to brain states

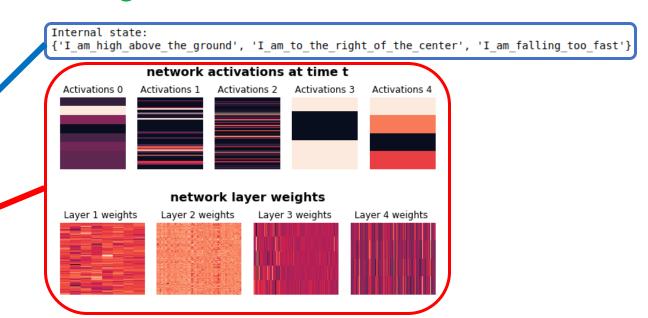
Phenomena of type A are ontologically reducible to phenomena of type B if and only if A's are nothing but B's

- V0
 - ✓ Internal states are causally reducible to brain states
 - Internal states are ontologically irreducible to brain states

Phenomena of type A are ontologically reducible to phenomena of type B if and only if A's are nothing but B's

Our ontology

- Layer weights of the neural network
- Connectivity of the neural network
- Activations of the neural network at time t
- The agent's observation at time t
- The agent's action at time t
- The position and velocity of the agent at time t
- A region the agent believes it's in
- Brain state at time t (set of layer weights, activations, and connectivity)
- Internal state at time t (set of regions the agent believes it's in)



- V0
 - ✓ Internal states are causally reducible to brain states
 - ✓ Internal states are ontologically irreducible to brain states

Phenomena of type A are ontologically reducible to phenomena of type B if and only if A's are nothing but B's

Our ontology Layer weights of the neural network Connectivity of the neural network Activations of the neural network Activations of the neural network The agent's observation at time The agent's action at time The position and velocity of the neural network Activations of the neural network Internal state instances are not "nothing but" brain state instances under our ontology (they are different classes) Brain state at time t (set of layer weights, activation), and connectivity)

• Internal state at time t (set of regions the agent believes it's in)

Is that the "real" ontology though?

- V0
 - ✓ Internal states are causally reducible to brain states
 - X Internal states are ontologically irreducible to brain states

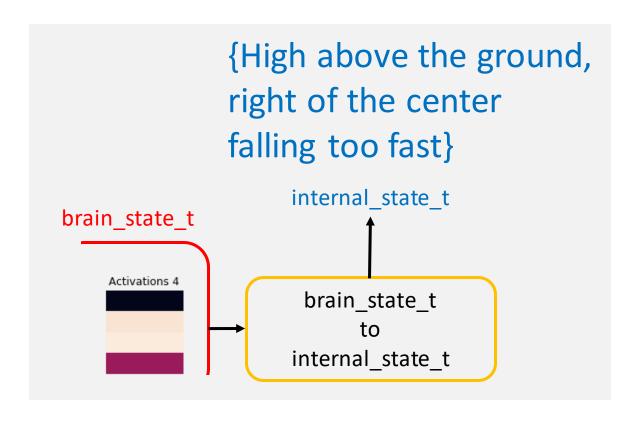
Phenomena of type A are ontologically reducible to phenomena of type B if and only if A's are nothing but B's

Our ontology

- Layer weights of the neural network
- Connectivity of the neural network
- Activations of the neural network at time t
- The agent's observation at time t
- The agent's action at time t
- The position and velocity of the agent at time t
- A region the agent believes it's in
- Brain state at time t (all of the bits contained in my computer)
- Internal state at time t (set of regions the agent believes it's in)

- Bits
- Python objects
- Electrons
- Quarks
- ..

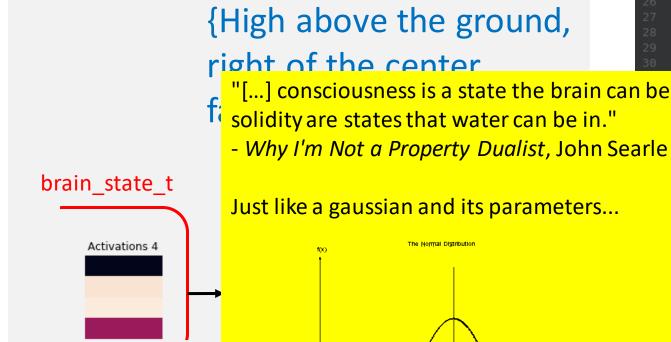
What's the deal with that function?



- Is this just some representation of "data flow"?
- Is this something closer to summarization?
- (or both?)

```
def brain state to internal state(brain state):
    def i am high above the ground(observation):
        return observation[1] > 0.5 # observation[1] accesses y position
   def i am low to the ground(observation):
        return observation[1] <= 0.5 # observation[1] accesses y position</pre>
   def i am to the right of the center(observation):
        return observation[0] > 0. # observation[0] accesses x position
   def i am to the left of the center(observation):
        return observation[0] <= 0. # observation[0] accesses x position</pre>
   def i am falling too fast(observation):
        return observation[3] < -0.2 # observation[0] accesses v velocity
    regions = [
        i am high above the ground,
        i am low to the ground,
        i am to the right of the center,
        i am to the left of the center,
        i am falling too fast
    internal state = set()
    recurrent activations = brain state['activations'][3]
    for activation, region in zip(recurrent activations, regions):
        if activation > 0.5:
            internal state.add(region. name )
    return internal state
```

What's the deal with that function?



- Is this just some
- Is this something
- (or both?)

$$\hat{\sigma}^2 = \frac{1}{n-1} \sum (X_i - \bar{X})^2$$

```
1  activation > 0.5:
       internal state.add(region. name )
return internal state
```

High above the ground,
$$\frac{26}{27}$$
 def i am low to the ground(observation): return observation[1] > 0.5 # observation[1] accesses y position def i am low to the ground(observation): return observation[1] > 0.5 # observation[1] accesses y position def i am low to the ground(observation): [1] accesses y position solidity are states that water can be in." - Why I'm Not a Property Dualist, John Searle

Just like a gaussian and its parameters...
$$\hat{\mu} = \bar{X} = \frac{1}{n} \sum X_i$$

5'][3]

ions, regions):

Conclusion

- Software engineer style philosophy reifying seemed to work well
- Created a V0 software agent who's
 - Internal states are causally reducible to brain states
 - Internal states are ontologically irreducible to brain states
- Download and play with the code yourself
 - https://github.com/Josh-Joseph/tsc-2019
- Disagree with us?
 - Great! Open an issue and/or submit a pull request in GitHub
- Thoughts on other theories of mind/consciousness that may be particularly well suited for this type of approach?