

# Finding causal theories quickly enough

## Building a responsive Apperception Engine

Jean-François Cloutier  
Robotics and Embodied project  
Active Inference Institute

Last updated: 2024-04-24

### Abstract

We describe the implementation of an Apperception Engine for use in a distributed cognitive architecture targeting Lego robots. An Apperception Engine applies symbolic machine learning to produce causal theories from short sequences of observations. The causal theories are small, human-readable logic programs that model an environment's latent generative processes and predict incoming observations. Finding these causal theories quickly is a difficult problem due to the size of the search space. We built an Apperception Engine that can infer good enough causal theories quickly enough to be useful to a robot actively engaging with its environment.

## 1. The need for an Apperception Engine

We implemented an Apperception Engine as part of an ongoing project to animate Lego robots with a distributed model of basal cognition. The principal goal of the project is to enable robots to learn, without human supervision, how to engage survivably with their physical environment.

The Apperception Engine, here described, is based on the original design from “[Making sense of sensory input](#)” [2] by Richard Evans et al who coined the name.

In our implementation of a cognitive architecture, we aim to conform to the principles of Active Inference but eschew a probabilistic implementation. Instead, we employ symbolic machine learning, specifically an Apperception Engine, to realize predictive processing capabilities.

An Apperception Engine is software that examines a short history of sensory inputs (observations) and discovers good enough causal theories that make sense of the observations by conforming to Immanuel Kant's conditions for [Unity of Apperception](#)[1]. A causal theory is a learned, human-readable logic program that can infer future observations.

The cognitive model making use of the Apperception Engine is a [Society of Mind](#) [3] in the sense that it is a collective of self-organizing and inter-dependent “cognition actors”. The cognition actors are processes under constant threat of extinction and thus motivated to collectively become competent at discovering and extracting the resources they need to

survive<sup>1</sup>. To that end, each cognition actor continually strives to make sense of what causes its sensory inputs (observations) and to decide how to act to bring about beneficial sensations/observations.

Each cognition actor is a separate process with its own limited scope of perception and enaction. It senses a small set of signals emitted by its environment which is either the outer world inhabited by the robot-as-autonomous-agent or the inner world populated by other cognition actors. A detailed description of our implementation of a society of mind will be the subject of a future publication<sup>2</sup>.

Each cognition actor needs to make sense of its observations (sensory inputs) in order to develop effective action policies to beneficially exploit and explore its environment. A cognition actor calls upon the Apperception Engine to devise or revise a causal theory that explains its more recent observations. Note that any action prompted by a cognition actor also becomes a signal it observes. This causal theory is enlisted to predict incoming signals and also to predict the consequences of potential actions.

To be usable, the Apperception Engine must be able to produce good enough results quickly, or else cognition actors might be constantly saddled with obsolete causal theories, possibly leading to their collective extinction. The Apperception Engine need not discover the best possible causal theories but it *must* operate in “useful time”, meaning the causal theories it looks for must still be relevant by the time they are found.

---

<sup>1</sup> A simulated metabolism is implemented in which computations progressively deplete a robot's initial and simulated energy stores. The robot can replenish its stores by positioning itself over different colored areas of the floor representing food (parking itself over one color is not enough). If energy stores become too low, cognition actors are pruned from the robot's society of mind to reduce the computational energy they draw. If and when all cognition actors are decimated, the robot effectively dies. The robot survives by learning, before it is too late, how to find the kinds of “food” it needs. Learning is done by evolving a society of mind that is, ideally, both effective and frugal at perceiving and engaging the robot's outer and inner world. This experimental setup will be detailed in a future publication.

<sup>2</sup> A quick overview: Cognition actors interact with each other by emitting and receiving predictions, prediction errors and action requests. Cognition actors also emit the beliefs they synthesize from their perceptions (the uncontested predictions they make plus the prediction errors they receive). These beliefs are signals that other cognition actors can sense and make predictions about so as to accumulate perceptions and then form their own, more abstract beliefs. Cognition actors thus form abstraction hierarchies.

Cognition actors also emit “cognitive signals” about fluctuations in their learning and sense making abilities. “Meta-cognition actors” can sense these signals and derive from them beliefs about the “cognitive aptitudes” of cognition actors. They can then act upon these beliefs to influence the cognitive functions of the individual cognition actors that they sense. Cognitive acts include the retiring of cognition actors and the instantiating of new ones. This cognition about cognition, we argue, realizes a form of *basal consciousness*. Meta-cognition actors themselves emit cognition signals that other meta-cognition actors can sense, thus forming “meta-cognition hierarchies”. Since all cognition actors can form abstraction hierarchies, meta-cognition actors can participate in both meta-cognition and abstraction hierarchies. This cognitive architecture is under active research. It is hoped that such a self-organizing society of mind will, through active engagement, be capable of acquiring behaviors of increasing competency needed to survive.

Our objective was to achieve good-enough response times so that individual cognition actors can use the Apperception Engine to make sense of their environment *as they interact with it*. Our Apperception Engine implementation was successfully tested on simple yet non-trivial apperception tasks.

## 2. What an Apperception Engine does

Succinctly put, an Apperception Engine is a program that takes a temporal sequence of observed states and searches for a logic program that, given one state, infers the next one. The logic program thus embodies a causal theory for the observations; it is a model of a latent generative process.

An observed state is a set of simultaneous observations about objects. An observation can be a property of an object or a relationship between two objects. States are snapshots that form a discrete (step-wise) temporal sequence.

To be considered a *valid* causal theory, the logic program must obey a number of unity constraints from Immanuel Kant's epistemology. These unity constraints ensure that the logic program "hangs together" and has explanatory power.

To help make our discussion concrete, we will be using the LEDs example from the Evans paper<sup>3</sup>.

Consider two visible LED lights. We'll call them A and B. An LED is either on or off whenever observed. In any given state taken from a sequence of observations, A and B will either be both on, both off, or one on and the other off. If we observe the blinking LEDs long enough, we notice that they go through a repeating blinking sequence.

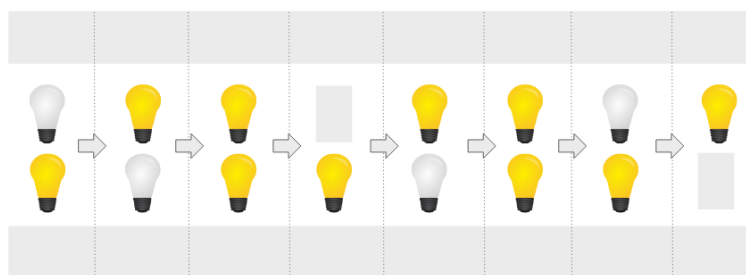


Figure 1 - Incomplete observations of two flashing lights

Observed states may be incomplete; data might be missing. The state of an LED might, once in a while, be unknown, or there might be other LEDs that are occluded and never observed, or there might be relationships between the LEDs that are also occluded (for example, they might be wired together out of sight.)

---

<sup>3</sup> In the paper, the LEDs are called "sensors". We find LEDs more illuminating.

The task of the Apperception Engine is to discover a causal theory that explains why the LEDs blink the way they do. A good theory will identify the LEDs in play (observed or not), the relevant properties of these LEDs (observed or not) and the relevant relationships between the LEDs (observed or not). It will also articulate rules in terms of these properties and relationships, rules that not only can predict incoming observations *but also* make sense of them<sup>4</sup>.

There are three types of rules in a causal theory produced by an Apperception Engine: Static rules, causal rules and conceptual constraints.

Static rules dictate what object properties and relationships are simultaneously allowed. An example might be that an LED must always be on if the one to its right is off. Static rules check if an inferred state is possible. They can also be used to infer missing observations in an incompletely observed state.

Causal rules dictate what the next observed state will be given a current state. An example might be that an LED will turn on if the one to its left is currently off. Causal rules are used to infer a next state from a current state. If applied to the last observations, they jointly predict incoming observations. If applied to a past observed state, they can backfill the missing state that follows it. So causal rules can both predict and retrodict.

Conceptual constraints are rules on relationships and properties. An example might be that any LED must have one and only one LED to its left, or that no LED can be simultaneously on and off<sup>5</sup>.

In addition to rules and constraints, a causal theory identifies initial conditions from which the rules are applied to create a trace. A trace is a sequence of states that, if the causal theory is 100% correct, recreates the observed states, possibly adding missing or unobserved relations and properties that themselves may involve unobserved objects.

A fully successful theory will create a trace that completely covers observed states. A not so successful theory will create a trace that only partially covers observed states.

---

<sup>4</sup> Predicting incoming observations is not, in and of itself, the same as making sense of them.

<sup>5</sup> In our implementation, mutual exclusion constraints on properties are implied.

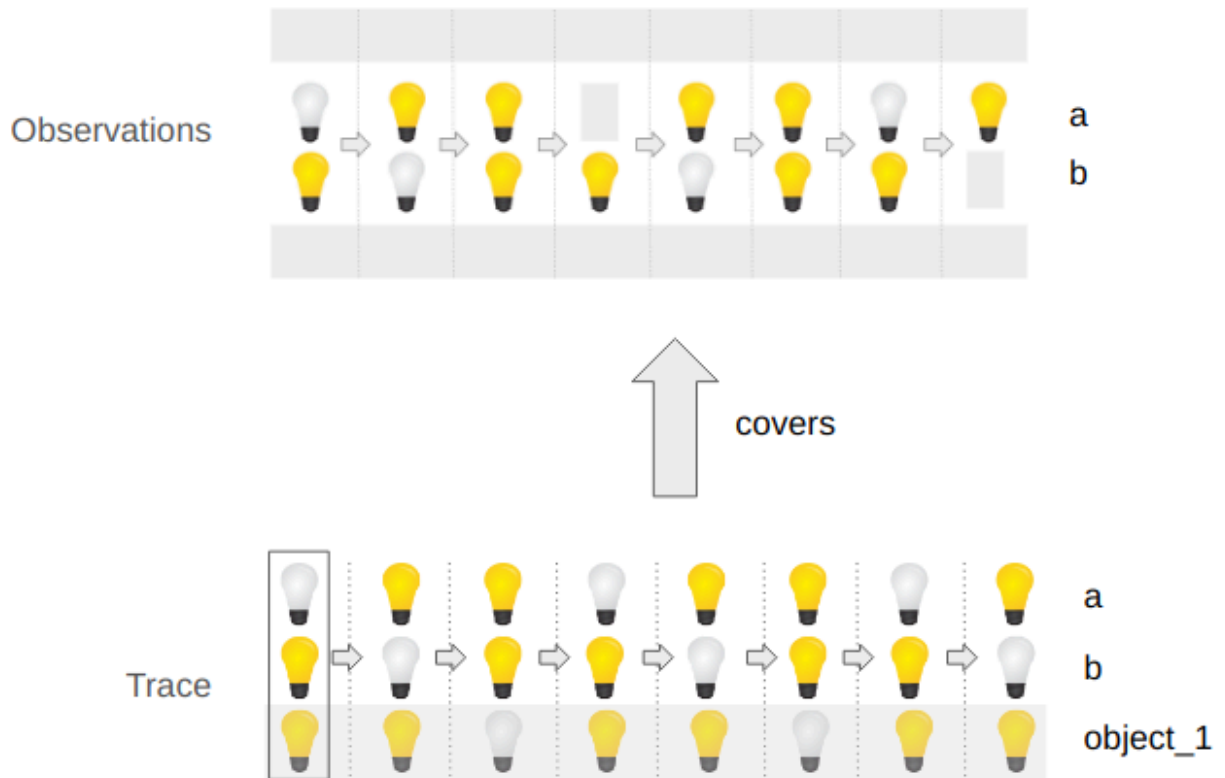


Figure 2- Trace covering observations of flashing lights

A causal theory is rated by how well it overlays the observed states and by how complex it is. The fewer and shorter the rules, the less complex the theory is.

The Apperception Engine is looking for a theory that maximizes coverage and minimizes complexity (Occam's razor).

The rules, conceptual constraints, observations and initial conditions of a causal theory are all expressed in terms of a vocabulary of

- object types (e.g. LED),
- objects (e.g. LED A, LED B, LED object\_1),
- property types (e.g. on(some LED, true or false),
- relation types (e.g. left-of(some LED, some other LED) and in-front-of(some LED, some other LED))

A "base" vocabulary is implied by the observations that are given as input to the Apperception Engine. For example, if all that is observed is the on status of two LEDs, then the two LEDs and the "on" property make up the base vocabulary.

The Apperception Engine may need to extend this base vocabulary. For example, it might want to add an unobserved LED (*object\_1* in the illustration above) and an unobserved relation in

order to craft a causal theory that successfully explains observations even though these observations are expressed in a more limited vocabulary.

Even if the Apperception Engine finds a causal theory with perfect coverage and minimal complexity, it will nevertheless reject it if the theory is not *unified*. Only a unified causal theory can be said to *make sense* of the observations.

A causal theory (initial conditions, static and causal rules, conceptual constraints) is unified if it meets spatial, conceptual, static and temporal unity conditions. The unity conditions are taken from Kant's Synthetic Unity of Apperception<sup>6</sup>.

**Conceptual unity:** Every type of predicates (predicates express properties and relations) must be referenced in at least one conceptual constraint. For example, a constraint on the relational "left of" predicate type might be that one and only one LED can be left of another. As for property types of predicates, they are always constrained in our implementation: An object's property can take only one value at once from an enumerated domain. For example, an LED's "on" property can be true or false but never both simultaneously<sup>7</sup>.

**Spatial unity:** Each object must be connected transitively to all others via relations. For example, if, in every state of the trace produced by the theory, the LEDs form a connected circle, then each one is related, directly or indirectly, to all others.

**Static unity:** Observations in a state (in a trace generated by a causal model or from observations) must obey all static rules. For example, if a static rule states that an LED must be on whenever the LED to its right is off, then LED A can not be off when LED B to its right is off.

**Temporal unity:** All states in a trace must be producible by first applying the causal rules on the prior state, and then carrying over what remains unchanged. For example, consider a state containing both on/off observations and left-of observations. If causal rules only change on/off states, the left-of observations would be immutable and thus would be carried unchanged throughout a temporal sequence of states.

To restate, the task of the Apperception Engine is, given a sequence of observations (possibly incomplete), to find a simple and unified causal theory composed of initial conditions, static rules, causal rules and conceptual constraints such that the causal theory-as-logic-program produces a trace from its initial conditions that covers as closely as possible the sequence of observations.

---

<sup>6</sup> See [The Apperception Engine](#) [4] for an in-depth discussion.

<sup>7</sup> The quantum realm is off-limits to the Apperception Engine.

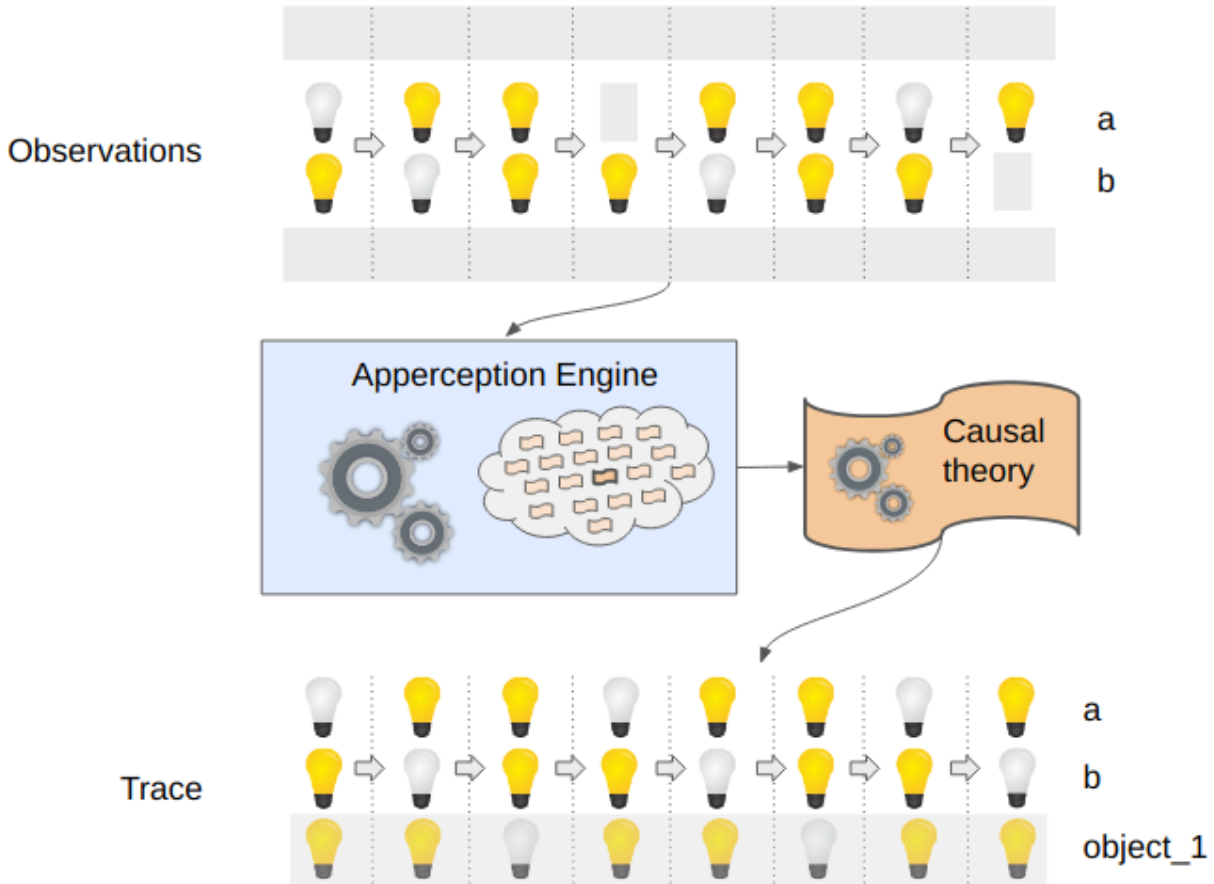


Figure 3 - Observations of flashing lights are given as input to the Apperception Engine tasked with discovering a causal theory that can produce a trace covering the observations.

Because the Apperception Engine does not restrict the causal theory strictly to observed objects, properties and relation, a successful theory might involve properties and relationships of hidden objects. For example, it is possible to set up a sequence of observations that only contain the on/off status of 2 LEDs such that articulating a successful causal theory will require the introduction of a hidden third LED plus a hidden relationship between the LEDs.

### 3. Requirements

While the functionality of our Apperception Engine is essentially the same as the Apperception Engine described in Evans et al., our performance requirements differ sharply.

The overall goal of the Evans implementation was to demonstrate that an inductive and symbolic approach to machine learning can outperform artificial neural networks on some classes of problems both in its economy (very little training data is required) and in the predictive capabilities demonstrated after the learning phase.

While some efforts were made to optimize their implementation, it was not a requirement that a causal theory be produced quickly. Finding a 100% successful causal theory with their Apperception Engine might require many hours of computation, perhaps even days<sup>8</sup>.

Our requirements are different. Cognition actors will call upon the Apperception Engine while engaging with their environment. They can not afford to wait long to replace obsolete causal theories with ones that better fit the latest observations. Our Apperception Engine needs to find good enough causal theories in useful time, which means seconds, not minutes, and certainly not hours or days.

It is important to note that a cognition actor will be restricted to a relatively small “umwelt”; the number of observed objects will be kept small and so will the variety of observed properties and relations. Furthermore, restrictions will be imposed on how many unobserved objects and object types, as well as types of properties and relationships, a cognition actor can imagine (abduce) when attempting to formulate a causal theory.

In short, our Apperception Engine is designed to find good enough solutions to relatively simple apperception problems in a matter of seconds.

## 4. Apperceiving in useful time is hard

Implementing a program that finds/builds logic programs that regenerates its input<sup>9</sup> is not, in and of itself, a trivial endeavor. But what makes implementing a functional Apperception Engine a truly hard problem is the sheer number of logic programs that are candidates for causal theories and that must be sifted through. This is the proverbial needle in a haystack.

---

<sup>8</sup> According to documentation in their [GitHub repository](#) [5].

<sup>9</sup> What an Apperception Engine does, in generic terms.



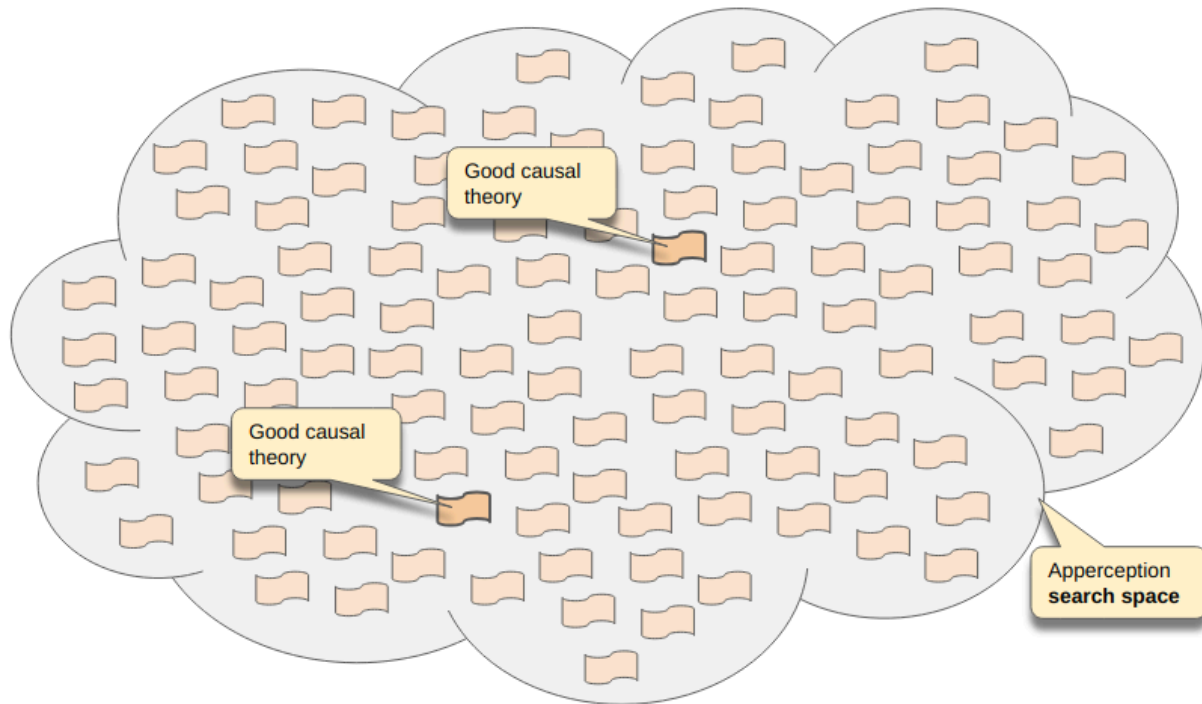


Figure 4- A search space for causal theories

If there was an algorithm that could, from a sequence of observations, directly compose, step by step, a causal theory, then the Apperception Engine could efficiently go directly from observations to causal theory. There is, to our knowledge, no such algorithm. What we are left with is a search problem.

A simplistic implementation would generate every potential causal theory in turn and test it for correctness (how well the trace it produces covers the sequence of observations), validity (unity), and for complexity (how many rules it has and how big they are.)

The problem is that “every potential causal theory” is an extremely large set. Let’s examine why.

A causal theory is composed of initial conditions (one potential observed state), a set of static rules, a set of causal rules, and a set of conceptual constraints.

There are typically many possible permutations of object properties and inter-object relations that can compose the initial conditions of a causal theory.

For every possible initial conditions, there are many ways in which the rules and constraints making up the rest of a causal theory can be composed.

Rules, whether static or causal, are of the form

```
head if body
```

where the head is either a property or a relation conditionally true of all objects, and the body is one or more properties/relations describing these objects and expressing the conditions.

For example,

```
on(led(X), true) if (left-of(led(X), led(Y)) and on(led(Y), false)
```

How many different individual rule formulations are possible given one vocabulary? A great many. How many combinations of such rules are possible? Also a great many.

There can be many formulations of conceptual constraints about the relation types in the vocabulary.

For example, every (LED) object has one and only one other (LED) object to its left:

```
one_related(left-of)
```

Or an (LED) object can be either left-of or in-front-of another, but not both:

```
one_relation([left-of, in-front-of])
```

A causal theory could have many different sets of conceptual constraints.

The number of possible causal theories given a vocabulary is roughly the number of possible initial conditions times the number of possible sets of static rules times the number of possible sets of causal rules times the number of possible sets of conceptual constraints.

Clearly, the space of all possible causal theories grows exponentially with the size of a vocabulary and quickly becomes too large to be enumerated and validated one by one in an acceptable amount of time.

This is bad enough for one vocabulary. The Apperception Engine searches for a variety of vocabularies, each one extending the base vocabulary (implied in the observations) by adding any number of “imagined” (unobserved) object types, objects and property and relation types.

The search space of all possible causal theories formulated using all possible vocabularies is, for all intents and purposes, infinite.

The generation of a candidate causal theory requires significant computational effort. So does the validation of a candidate theory; a trace must be produced by applying the rules on the initial conditions and then the trace must be checked against the input sequence of observation. And the causal theory, as well as the trace it produces, must be checked for unity (spatial, conceptual and static and temporal.)

A simple “generate then validate” implementation would clearly take forever to systematically find and evaluate all candidate causal theories. Consequently, the odds of finding a valid theory early in the search process would be very small indeed.

One strategy employed by the Evans et al. implementation and also used in our own is to limit and partition the search space, and then restrict the search to partitions that would produce reasonably small theories.

The Apperception Engine partitions the search space into regions<sup>10</sup>, where each region imposes lesser or greater limits on how “imaginative” a causal theory can be, namely how many unobserved objects, properties and relations can be added to the vocabulary implicit in the observations.

The Apperception Engine does not explore all possible regions into which the search space is partitioned, only those that require not too much imagination (i.e. minimal extensions to the base vocabulary). For example, in the LEDs example, imagining one unobserved LED is sensible but imagining ten probably is not.

There can be many alternate vocabularies that fit within the “limits on imagination” set by a region of the search space, and so each region is subdivided into “causal theory templates”. Each template has its own vocabulary, extended from the base vocabulary, but within the limits set by its region.

A template also sets limits on the number and size of static and causal rules. These limits are commensurate with the vocabulary’s richness.

Each template is searched for successful causal theories expressed in the template’s vocabulary and that do not exceed a maximum complexity.

---

<sup>10</sup> Regions are called “tuples” in Evans et al.

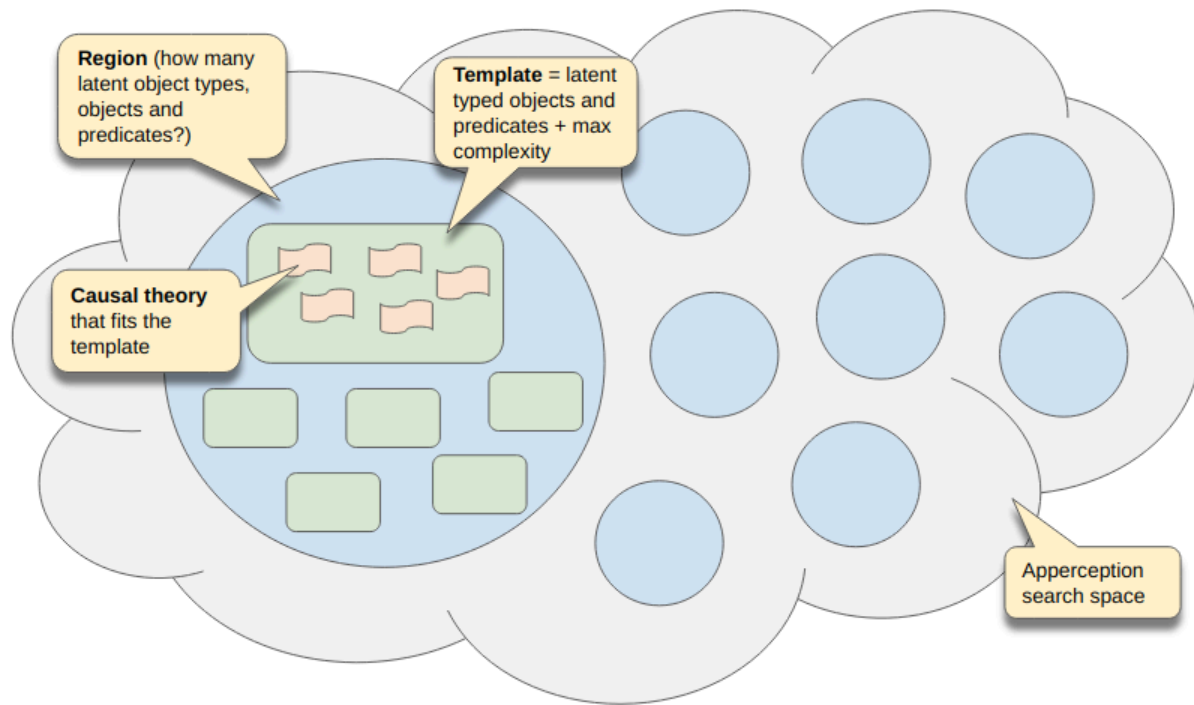


Figure 5 - Partitioning the search space into regions containing templates

The benefit of this partitioning and delimiting of the search space is that searches can terminate. The downside is that a satisfactory causal theory might not be found in the time allowed and within the selected regions of an otherwise infinite search space.

Partitioning and restricting the search space helps but it is not enough. We needed to do a lot more if our implementation of the Apperception Engine was to have a reasonable chance of finding good enough causal theories within seconds.

## 5. Implementation

Evans et al.'s implementation uses a variant of Datalog (a simple logic programming language) together with a constraint solver library written in Haskell. The emphasis appears to have been on keeping the code as simple as possible, which is desirable from a research perspective when performance is not a core concern.

In our [implementation](#) [6] performance is a central concern.

We implemented our Apperception Engine *de novo* in [SWI-Prolog](#) [7], augmented by a [Constraint Handling Rules](#) [8] (CHR) library.

[Prolog](#) [9] is a *logic programming* language. A Prolog program is a set of rules specifying how some facts can be proven from other facts. Running a program consists in asking it to prove an

implied fact, for example that there exists some causal theory that explains observations. Here, the observations are known but the causal theory is unknown and to be inferred.

A proof may succeed with a solution that instantiates the unknowns (e.g. a successful causal theory), or it may fail. The proof process is recursive; proving a fact may entail proving other facts that themselves require the proving of further facts etc., until no unknown remains.

During the proof process, a fact might be provable in different ways. Prolog will choose one way. If this leads to failure downstream, the proof process will backtrack and try another way. If failures persist and Prolog runs out of ways to prove the problematic fact, it backtracks to an earlier choice point and tries a different path forward from there. If it backtracks all the way to the start of the proof then the proof has failed and there is no solution.

Backtracking on failure is thus a way for a program to automatically try different paths to infer a solution until it is found or all possible paths are exhausted. Thanks to this capability, Prolog programs that traverse complex search spaces can be written in relatively few lines of code without sacrificing efficiency.

The Prolog programming language, with its deductive computing model and its use of backtracking, is particularly well-suited for implementing an Apperception Engine.

SWI-Prolog is an open-source Prolog with support for concurrency on multicore hardware. It enables our Apperception Engine to explore in parallel as many templates as the CPU has cores. This greatly accelerates the search.

Our implementation uses a combination of Prolog and Constraint Handling Rules. CHR is a “committed-choice rule-based language embedded in Prolog”. In effect, it adds abductive reasoning to Prolog’s deductive reasoning. It enables a program, here the Apperception Engine, to make assumptions and hold on to them unless and until further assumptions cause CHR rules to refute them.

For example, the Apperception Engine will make the assumption that a causal theory under construction is unified. The assumption holds until, say, a rule is added to the causal theory that contradicts a condition of unification (written in CHR).

An assumption, once contradicted, causes a failure that forces the program to backtrack to a point in its execution where the assumption still holds. The program can then resume from that point and go on to, for example, add a different rule that hopefully will not break the unity assumption.

This use of abductive reasoning, combined with backtracking on failure, enables the program to fail as soon as possible, pruning searches early, thus avoiding wasting time following a path doomed to fail to its dead end.

Our implementation of the Apperception Engine entangles three kinds of logical inference.

Prolog programs solve problems by logically *proving* that there are solutions. They do so by applying *deductive reasoning*, from known or inferable facts, using the logic rules that compose these programs. CHR adds *abductive reasoning* to Prolog; it allows programs to assume, during their execution, that some facts are true *without proof* until ruled otherwise.

The Apperception Engine, written in Prolog + CHR thus uses deductive and abductive reasoning to do what it does. And what it does is *inductive reasoning*; it builds causal theories out of observations (inferring rules from facts).

To more fully appreciate this “entanglement of logics”, consider that the causal theories built using these three kinds of logic are themselves logic programs that can *deduce* incoming observations from prior observations.

## 6. Tactics and heuristics

An Apperception Engine carries out “apperception tasks”. An apperception task is expressed as a sequence of observations for which a valid, unified causal theory is to be found. Finding a causal theory (proving that one exists) means searching a vast space of candidate causal theories for one that “works” and is unified, and thus makes sense of the observations given as input.

As mentioned earlier, a naive implementation, one that systematically generates each possible candidate causal theory and then tests it for accuracy and unity, would be grossly inefficient.

In order to find good enough causal theories in useful time (in seconds), we needed to employ a number of tactics and heuristics. The common objective of all tactics and heuristics is to increase the odds that exploration of the search space will yield an acceptable solution in a small amount of time.

### 6.1 Fail early

If a path taken through the search space is going to fail, it is best to figure this out as soon as possible and abandon that path before any more time and computing resources are wasted.

Our implementation makes use of CHR to check that validity assumptions are contradicted as soon as possible in the process of building (i.e. searching for) a causal theory.

For example, one validity constraint is that a static rule must not be recursive because that makes it tautological. A recursive rule would say something like “An **LED X** must be **on** if there is an LED Y to its left and **LED X** is **on**”. This failure condition is checked while each rule is being built and can be triggered even before a rule is done being built. The implementation does

not wait for all of a theory's static and causal rules to be built, or even a single one, to check for recursive static rules, because a lot of work would then need to be undone.

## 6.2 Lower the bar

Looking for the perfect causal theory can be prohibitively expensive because they are rare whereas good enough theories might be plentiful and thus much faster to find.

Our implementation allows apperception parameters to specify what a good enough causal theory is. For example, an apperception task might be launched such that it is good enough for a causal theory to have a coverage of 80%, meaning the trace it produces covers 80% of the sequence of observations given as input.

As soon as a good enough causal theory is found in some template, the overall search is stopped and the successful causal theory is produced.

By lowering the bar for success, the search can terminate quickly with a useful, if imperfect, causal theory.

## 6.3 Avoid repetition

When traversing a search space, there is a lot of potential for duplication of work. Retracing one's steps is a surefire way of wasting time and resources.

For example, in a causal theory, the static and causal rules are restricted to conjunctions and are of the form "A if B and C and D". If care is not taken, another path through the search space might try the alternate rule "A if D and C and B". Both rules are logically equivalent so only one variant needs to be tried.

Our implementation tries to avoid repetitive work.

## 6.4 Throw the dice

There are a lot of choice points when building/searching for a causal theory. Which property or relation to use in the head of a rule or next in constructing the body of the rule? How many predicates should appear in the rule? How many static rules should the causal theory have? Etc. When faced with a choice, instead of trying options in the same order every time, our implementation randomizes the order in which different options are tried.

Consequently, running the same apperception task multiple times does not yield the same causal theories in the same order. Sometimes a particularly apt causal theory will be found early in the search, sometimes late. This is arguably better than possibly always late.

## 6.5 Set time limits

The Apperception Engine divides the search space into regions (each region sets limits on imagination) and those into templates (extended vocabularies and limits on complexity). A lot of time could be spent searching in one template, leaving little to no time to explore others. To avoid this, our implementation (and Evans et al.'s as well) time boxes searching each template<sup>11</sup>.

In our implementation, the time limit is a function of the size of the vocabulary attached to a template.

When the time allocated to a template is exhausted, the best causal theories found in it are gathered to compete against the best theories found in other templates.

An apperception task is also given an overall time limit. Once that time limit is reached without finding a good enough causal theory, the overall search is terminated and the runner-up theories found in the explored search space are produced as best-effort solutions.

## 6.6 Be selective

A number of heuristics fall under this category. The aim is not to look everywhere in a sprawling search space, but instead to look where, if a solution could be found, it might be found quickly.

Our implementation, and Evans', favors simple causal theories over more complex ones. For one thing, they are cheaper to build and evaluate. Since the Apperception Engine favors **simplicity** (Occam's razor), it is best to search for simple theories first. This means favoring templates with small vocabularies and with severe restrictions on the allowed number and sizes of rules composing the causal theories.

Furthermore, if a search in a template comes up empty after a short amount of time, our implementation presumes that spending more time in it is a fool's errand, becomes **impatient**, gives up and moves on to another template.

Our implementation does not try every possibility when building/searching for a causal theory; it instead **samples** only a few.

For example, when assembling initial conditions (from which a trace will be built by applying the theory's rules), there are a large number of permutations of initial observations to choose from. Our implementation tries only a few under the assumption that they are likely to all produce equivalent traces.

Another example of sampling is **limiting the breadth of search** when building a set of static or causal rules. For each candidate first rule in a set of rules, there can be a very large number of

---

<sup>11</sup> Our time limits are expressed in seconds, not hours.



candidates for the second rule and then third rule etc. If the breadth of search for the second, third etc. rule is not limited, much of the time allocated to searching a template might be used exploring “under” the first candidate rule looked at. The Apperception Engine limits the search for the next rule in a set to, say, 100 candidates, so as to leave time for searching under alternate first, second etc. rules.

Our implementation uses **iterative deepening**, which is another way of selectively searching for a satisfactory causal theory. If the overall time budget is not exhausted after the initial search across templates, the more promising templates that have not been fully explored are retained, and are searched again, this time with more time allocated to each one. The hope is that incompletely searched templates that produce better causal theories might produce even better ones if given more time. And because luck is involved (see “Throw the dice”), a new iteration could get lucky and find, in a promising template, a satisfactory theory that eluded prior iterations.

## 7. Apperception parameters

An apperception task is accompanied by parameters that limit the search in terms of its scope, success criteria and time budget.

The parameters are:

- How much time is allocated to the apperception task overall.
- The minimum coverage that defines a satisfactory causal theory - the search stops when the first satisfactory theory is found.
- How many runner-ups causal theories should be produced if no satisfactory theory is found in the time allowed.
- What are the overall limits imposed on imagination (up to how many unobserved objects, object types, relation types and property types can the Apperception Engine invent to express a causal theory.)
- How quickly should a new iteration narrow down to the most promising templates when a search completes without success and there’s still time left.

These parameters impact the likelihood of quickly finding a satisfactory causal theory. They can be tuned for different apperception tasks.

## 8. Sample apperception tasks

We present the results of running two apperception tasks. The first task is to explain the blinking pattern of two LED lights (the LED task). The other task is to find the rule(s) of an Elementary Cellular Automaton (the ECA task).

Each task defines search spaces with significantly different properties. The LED task involves few objects (two lights with perhaps another one that’s unobserved), possibly connected by unobserved relations. The ECA task involves many more objects (11 cells, either on or off, all

observed) connected by an observed relation (every cell has exactly one cell to its left, effectively forming a circle.)

## 8.1 LED task

The Apperception Engine is given a sequence of observations expressed as predicates.



```
sensed(on, [object(led, a), false], 1).
sensed(on, [object(led, b), true], 1).
sensed(on, [object(led, a), true], 2).
sensed(on, [object(led, b), false], 2).
sensed(on, [object(led, a), true], 3).
sensed(on, [object(led, b), true], 3).
sensed(on, [object(led, b), true], 4).
sensed(on, [object(led, a), true], 5).
sensed(on, [object(led, b), false], 5).
sensed(on, [object(led, a), true], 6).
sensed(on, [object(led, b), true], 6).
sensed(on, [object(led, a), false], 7).
sensed(on, [object(led, b), true], 7).
sensed(on, [object(led, a), true], 8).
```

It was then asked to prove that a causal theory exists with these parameters:

```
MaxSignatureExtension = max_extension{max_object_types:1, max_objects:1,
max_predicate_types:1},
ApperceptionLimits = apperception_limits{max_signature_extension:
MaxSignatureExtension, good_enough_coverage: 100, keep_n_theories: 3, funnel:
3-2, time_secs: 60},
```

With `max_extension{max_object_types:1, max_objects:1, max_predicate_types:1}`, the Apperception Engine is allowed to imagine up to one hidden object type, object and relation or property.

With `good_enough_coverage: 100`, it looks for a fully accurate causal theory. If not found, it will return the 3 runner ups (`keep_n_theories: 3`). The search is given a maximum of 60 seconds to find a satisfactory causal theory (`time_secs: 60`).

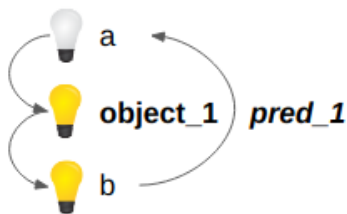
This is the successful causal theory the Apperception Engine typically finds first:

```

static_constraints:[one_related(pred_1)]
static_rules:[on(_A, true)-[pred_1(_B, _A), on(_B, false)]]
causal_rules:[on(_C, false)-[pred_1(_C, _D), on(_D, false)]]
initial_conditions:[pred_1(object_1, b), pred_1(b, a), pred_1(a, object_1), on(object_1, true), on(b, true)
on(a, false)]

```

Interpreting the output a bit demonstrates that the causal theory does provide a cogent explanation of the blinking sequence given as input.



A light is related to one and only one other light via *pred\_1*

A light is on whenever another that is *pred\_1* to it is off

A light turns off if it is *pred\_1* to another light that was off

Figure 6 - A causal theory making sense of observed flashing lights

In order to make sense of the observations, the Apperception Engine imagined a hidden LED (it named it *object\_1*) as well as an unobserved relationship (it named it *pred\_1*) that connects all LEDs together.

The causal theory is unified spatially, conceptually, statically and temporally.

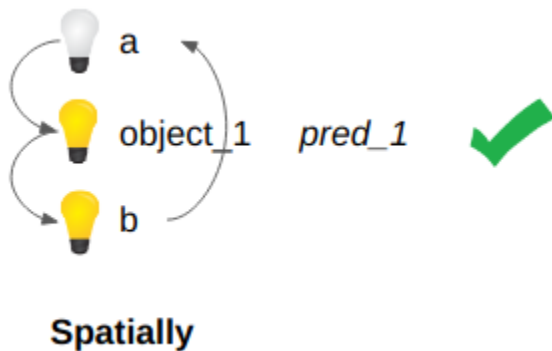


Figure 7- The causal theory is spatially unified

All objects are connected directly or transitively via the *pred\_1* relation



### Conceptually

Figure 8 - The causal theory is conceptually unified

The *pred\_1* relation is constrained in that only one LED can be *pred\_1* to another (the *one\_related* conceptual constraint). The *on* property is implicitly constrained in our implementation in that an LED's *on* property can only take one value (true or false) at a time.

**Statically** A light is on whenever another that is *pred\_1* to it is off

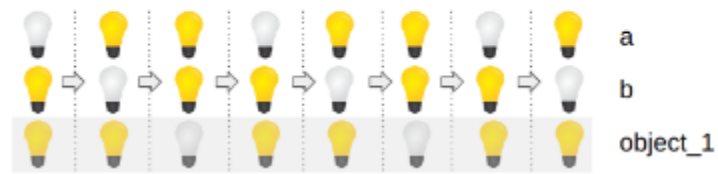
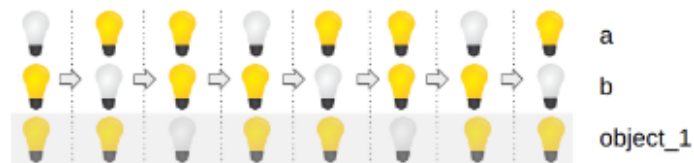


Figure 9 - The causal theory is statically unified

The static rule is respected in every traced state.



**Temporally** A light turns off if it is *pred\_1* to another light that was off

Figure 10 - The causal theory is temporally unified

The causal rule describes correctly what makes LEDs turn on or off in the trace.

Twenty consecutive searches were run.

## Time elapsed (seconds)

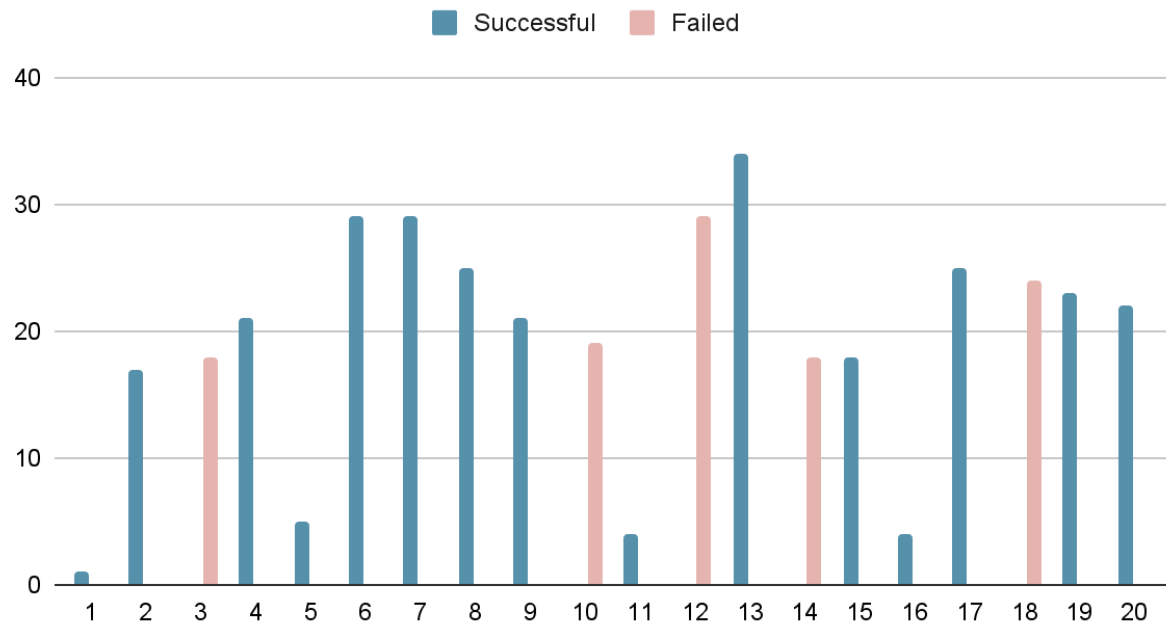


Figure 11 - Results from 20 consecutive searches for an accurate causal theory making sense of observed flashing lights

A 100% coverage causal theory was found in 15 out of 20 runs. When the run was successful, it took between 1 and 34 seconds to find a solution. The average time elapsed was 18 seconds with a standard deviation of  $10^{12}$ .

Five runs failed to find an accurate theory; they only found runner-up theories with 75% coverage<sup>13</sup>. These however were found after under 1 second, though it took on average 21 seconds for the Apperception Engine to give up the search for an accurate causal theory.

## 8.2 ECA task

The task is to uncover the rule (or rules) that drive an Elementary Cellular Automaton<sup>14</sup> (ECA).

*“an **elementary cellular automaton** is a one-dimensional cellular automaton where there are two possible states (labeled 0 and 1) and the rule to determine the state of a cell in the next generation depends only on the current state of the cell and its two immediate neighbors”*

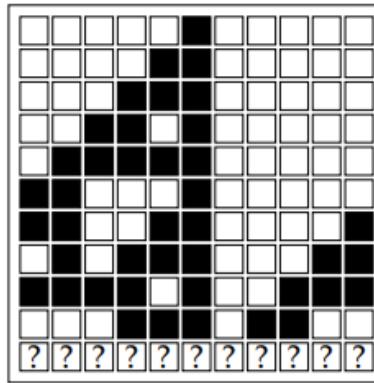
<sup>12</sup> The searches ran on a standalone computer with an 11th Gen Intel® Core™ i7-11800H @ 2.30GHz × 16.

<sup>13</sup> The Apperception Engine does not explore the entire search space. As a consequence, there is always a chance that a satisfactory solution will be missed.

<sup>14</sup> [https://en.wikipedia.org/wiki/Elementary\\_cellular\\_automaton](https://en.wikipedia.org/wiki/Elementary_cellular_automaton)

While an ECA is depicted as rows where each row succeeds the other in time, they are in fact circles with the first cell considered to be to the right of the last cell, that is, cell 2 is to the right of cell 1, ..., cell 11 is to the right of cell 10, and cell 1 is to the right of cell 11.

The Apperception Engine is given these observed rows of an ECA as input:



<https://arxiv.org/pdf/1910.02227.pdf> (p.34)

Figure 12 - Observed 10 consecutive states of an Elementary Cellular Automaton

Nothing is left to the imagination in this apperception task: All objects are observed (11 cells), their on/off states (black is on) are all accounted for in the observations, the only relation is *right\_of* and all the *right\_of* relations between cells are observed.

What is asked of the Apperception Engine, this time, is to find not a perfect solution (it would take much too long) but a solution that is at least 90% correct.

The key parameters given to the apperception task are

`max_extension{max_object_types:0, max_objects:0, max_predicate_types:0}` (nothing is left to the imagination), `good_enough_coverage: 90` (90% accuracy is acceptable), and `time_secs: 30` (time allocated is 30 seconds).

The Apperception Engine finds a causal theory that predicts the behavior of the ECA with 93% accuracy. It contains one causal rule and no static rules.

The causal rule is `on(_A, true) - [right_of(_B, _A), on(_B, true)]` which translates to “A cell will turn on if the cell to its right is on”<sup>15</sup>.

Ten consecutive searches were run. A causal theory with 93% accuracy was found reliably after 1 to 3 seconds of searching.

<sup>15</sup> The two causal rules of an 100% accurate theory are “A cell will turn on if it is off and the cell to its right is on” and “A cell will turn off if it is on and the cells to its right is on and left are on”

The trace produced by this good enough causal theory is

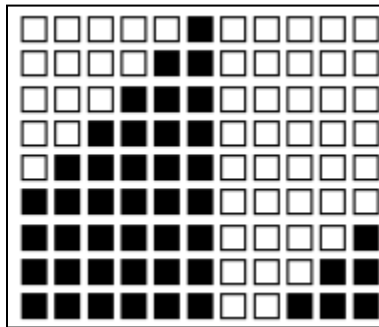


Figure 13 - Trace produced by an inaccurate but good-enough causal theory

It approximates the observations given as input while not quite reproducing them. However, the causal theory was found consistently in a handful of seconds.

When given enough time, the Apperception Engine can eventually find the simplest possible solution with 100% coverage, that is, one that produces a trace that perfectly recreates the observations. This ideal causal theory would contain one conceptual constraint, two causal rules and no static rule:

```
static_constraints: [one_related(right_of)]
static_rules: []
causal_rules: [
  on(X,true)-[right_of(Y,X),on(Y,true),on(X,false)],
  on(Y,false)-[right_of(Z,Y),right_of(Y,X),on(Y,true),on(X,true),on(Z,true)]
]
```

The causal theory says:

1. Exactly one cell is to the right of each cell
2. A cell that's off will turn on if the cell to its right is on
3. A cell that's on will turn off if both neighboring cells are on

The theory has no static rule; each new state is fully described by the on/off property changes the causal rules prescribe, complemented by the unchanged cells on/off properties and by the immutable relative positions of cells.

This causal theory is the well-known [Rule 110](#) [10] “which is neither completely stable nor completely chaotic.”

The import to the Apperception Engine is that a trace produced by the causal theory will likely not repeat for quite a while. Calculating the coverage of the observations by a candidate causal theory will thus be misleading unless the initial conditions from which the trace is produced are identical to the first state in the observations. This makes it much harder for the Apperception Engine to find a satisfactory solution.

Nonetheless, in one run<sup>16</sup> the Apperception Engine did find an accurate causal theory that is equivalent to the above except that it contains a superfluous condition in one of its causal rules.

```
static_constraints: [one_related(right_of)]
static_rules: []
causal_rules: [
  on(X,true)-[right_of(Z,Y), right_of(Y,X),on(Y,true),on(X,false)],
  on(Y,false)-[right_of(Z,Y),right_of(Y,X),on(Y,true),on(X,true),on(Z,true)]
]
```

The condition stating “there’s a cell to the right of the one that causes a cell to turn on” is correct but it is always true and is thus redundant. The rules, while not minimal, are nevertheless valid (they are unified). Even though the theory the Apperception Engine found is accurate, it is inferior to the ideal causal theory shown earlier because it is more complex.

## 9. Conclusion

We implemented an Apperception Engine based on the one described in [Evans et al. \[2\]](#) but with an emphasis on finding good enough causal theories quickly.

A causal theory is a small, human-readable logic program that can recreate a short sequence of observations and can thus predict future observations. Causal theories are symbolic, non-probabilistic generative models.

The intent is for cognition actors to call upon the Apperception Engine to make sense of their observations in “useful time” and thus be able to predict incoming observations given planned actions.

We achieved promising results.

---

<sup>16</sup> It found the causal theory after spending 230 seconds searching the template that contains it.



Further optimization might be needed as we move to implementing cognition actors in Prolog + CHR and have them use the Apperception Engine to make sense of their environment as they engage with it.

## 10. Closing thought

In [Kant's View of the Mind and Consciousness of Self](#) [1], we read "Translated into contemporary terms, the core of [Kant's transcendental method of inquiry] is *inference to the best explanation*, the method of *postulating unobservable mental mechanisms* in order to *explain observed behaviour*.<sup>17</sup>"

If we replace “unobservable mental mechanisms” by “latent processes”, we get an abstract description of how the Apperception Engine operates.

Thus what the Apperception Engine does - inducing causal theories of hidden processes from observations - is also how Kant derived the epistemology that's at the core of the Apperception Engine. Douglas Hofstadter might call this a “[strange loop](#)” [11].

---

<sup>17</sup> Emphasis mine

## Acknowledgements

We are very grateful to Daniel Friedman of the Active Inference Institute for the constant encouragement and for the fruitful discussions that helped propel this work.

## References

- [1] Kant's View of the Mind and Consciousness of Self, Stanford Encyclopedia of Philosophy, <https://plato.stanford.edu/entries/kant-mind/>
- [2] Richard Evans, Jose Hernandez-Orallo, Johannes Welbl, Pushmeet Kohli, Marek Sergot, Making sense of sensory input, arXiv preprint arXiv:1910.02227
- [3] Minsky, Marvin (1986). The Society of Mind. New York: Simon & Schuster. ISBN 0-671-60740-5
- [4] Evans, Richard (2022). The Apperception Engine. In Hyeongjoo Kim & Dieter Schöneck (eds.), Kant and Artificial Intelligence. De Gruyter. pp. 39-104.
- [5] RichardEvans/apperception GitHub repository, <https://github.com/RichardEvans/apperception>
- [6] jfcloutier/andy\_karma GitHub repository, [https://github.com/jfcloutier/andy\\_karma/tree/main/code/apperception](https://github.com/jfcloutier/andy_karma/tree/main/code/apperception)
- [7] SWI-Prolog <https://www.swi-prolog.org/>
- [8] Thom Frühwirth (2009), Constraint Handling Rules, Cambridge University Press, ISBN 9780521877763,
- [9] Pierre Deransart, AbdelAli Ed-Dbali, Laurent Cervon (1996), Prolog: The Standard: Reference Manual, Springer, ISBN 978-3764352998
- [10] Rule 110, Wikipedia, [https://en.wikipedia.org/wiki/Rule\\_110](https://en.wikipedia.org/wiki/Rule_110)
- [11] Hofstadter, Douglas R. (2013) [2007], I Am a Strange Loop, Basic Books, ISBN 978-0465008377