

Wandering in the Labyrinth of Thinking

– a minimalist cognitive architecture combining
reinforcement learning, deep learning, and logic structure

甄景贤 (King-Yin Yan) Dmitri Tkatch 肖达 (Da Xiao) Juan Carlos Kuri Pinto

General.Intelligence@Gmail.com

Abstract. The bottleneck algorithm in AGI is the inductive learning of knowledge. The importance of this algorithm in our AI era is like what the steam engine was to the industrial revolution. This paper: 1) outlines a minimalist AGI architecture; 2) defines a logic structure for AGI systems; 3) imposes the logic structure on the control system, as inductive bias to speed up learning.

Keywords: cognitive architecture, reinforcement learning, deep learning, logic-based artificial intelligence

0 Summary

We propose an AGI architecture:

1. with **reinforcement learning** (RL) as top-level framework
 - State space = mental space
2. **Logic** structure is imposed on the **knowledge representation** (KR)
 - State transitions are given by logic rules = actions in RL
 - The logic state x is decomposable into **propositions** (§2.4)
3. The set of logic rules is approximated by a deep neural network
 - Just the most basic kind of feed-forward neural network (FFNN) is required
 - Logic conjunctions are **commutative**, so working-memory elements can be presented in any order (§2.5)
 - **Stochastic** actions are represented by **Gaussian kernels** (radial basis functions) (§2.6), thus partly avoiding the curse of dimensionality

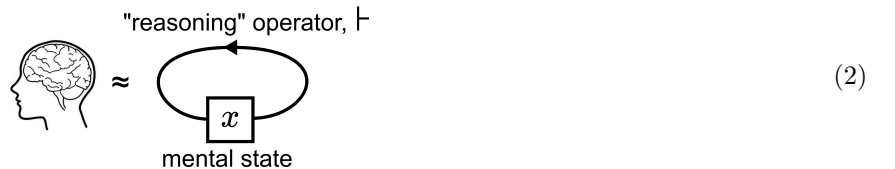
The rest of this paper will explain these design features in detail.

1 Reinforcement-learning architecture

The **metaphor** in the title of this paper is that of RL controlling an autonomous agent to navigate the maze of “thoughts space”, seeking the optimal path:



The main idea is to regard “thinking” as a **dynamical system** operating on **mental states**:



A mental state is a **set of propositions**, for example:

- I am in my room, writing a paper for AGI-2019.
- I am in the midst of writing the sentence, “I am in my room, ...”
- I am about to write a gerund phrase “writing a paper...”

Thinking is the process of **transitioning** from one mental state to another. As I am writing now, I use my mental states to keep track of where I am at within the sentence’s syntax, so that I can construct my sentence grammatically.

1.1 Actions = cognitive state-transitions = “thinking”

Our system consists of two main algorithms:

1. Learning the transition function \vdash or $F : x \mapsto x'$. F represents the **knowledge** that constrains thinking. In other words, the learning of F is the learning of “static” knowledge.
2. Transitioning from x to x' . This corresponds to “thinking” under the guidance of the static knowledge F .

In our architecture, F can be implemented as a simple feed-forward neural network (where “deep” simply means “many layers”):



Since a recurrent NN is Turing-complete, this can be viewed as a minimalist AGI. But its learning may be too slow without further **inductive bias** (*cf* the “no free lunch” theorem [18]) — so we will further modify F by imposing the logic structure of reasoning on it (§2.4 and §2.5).

In principle, every state is potentially **reachable** from every other state, if a logic rule exists between them. Now we use a deep FFNN to represent the set of all logic rules. This is a key efficiency-boosting step, because deep neural networks allows to use a polynomial number of parameters to represent an exponential number of mappings.

Note that parts of the state x would be reserved and directly connect to the **input** and **output** of the AGI system.


1.2 Comparison with AIXI [17]

AIXI’s environmental setting is the same as ours, but its agent’s internal model is a universal Turing machine, and the optimal action is chosen by maximizing potential rewards over all programs of the UTM. In our (minimal) model, the UTM is constrained to be a neural network, where the NN’s **state** is analogous to the UTM’s **tape**, and the optimal weights (program) are found via Bellman optimality.

2 Logic structure

2.1 Logic is needed as an inductive bias

We know that the transition function F is analogous to \vdash , the logic consequence or entailment operator. So we want to impose this logic structure on F .

By logic structure we mean that F would act like a **knowledge base**  containing a large number of logic **rules**, as in the setting of classical logic-based AI.

2.2 Geometry induced by logic rules

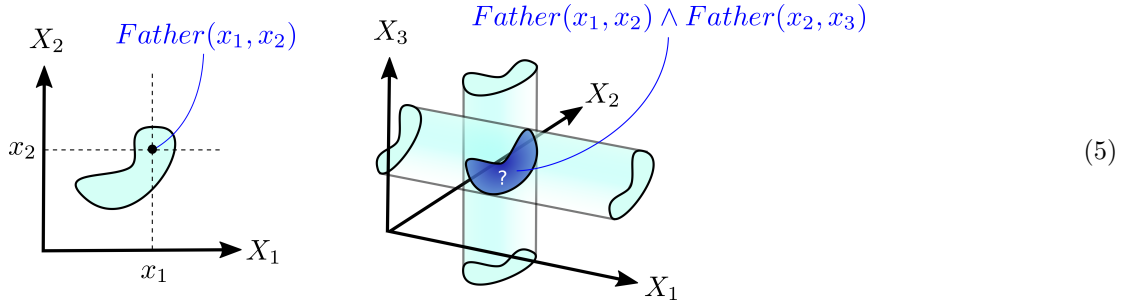
A logic rule is a conditional formula with variables. For example:

$$\forall X \forall Y \forall Z. \text{ father}(X, Y) \wedge \text{ father}(Y, Z) \Rightarrow \text{ grandfather}(X, Z) \quad (4)$$

where the red lines show what I call “linkages” between different appearances of the same variables.

Quantification of logic variables, with their linkages, result in **cylindrical** and **diagonal** structures when the logic is interpreted *geometrically*. This is the reason why Tarski discovered the **cylindric algebra** structure of first-order

predicate logic [14] [15] [1] [5] [6]. Cylindrical shapes can arise from quantification as illustrated below:



And “linkages” cause the graph of the \vdash map to *pass through* diagonal lines such as follows:



We are trying to use neural networks to approximate such functions (*ie*, these geometric shapes). One can visualize, as the shape of neural decision-boundaries approximate such diagonals, the matching of first-order objects gradually go from partial to fully-quantified \forall and \exists . This may be even better than if we fix the logic to have exact quantifications, as quantified rules can be learned gradually. There is also *empirical* evidence that NNs can well-approximate logical maps, because the *symbolic* matching and substitution of logic variables is very similar to what occurs in *machine translation* between natural languages; In recent years, deep learning is fairly successful at the latter task.

2.3 Form of a logic rule

So what exactly is the logic structure? Recall that inside our RL model:

- state $\mathbf{x} \in \mathbb{X}$ = mental state = set of logic propositions $P_i \in \mathbb{P}$
- environment = state space \mathbb{X} = mental space
- actions $\mathbf{a} \in \mathbb{A}$ = logic rules

For our current prototype system, an action = a logic **rule** is of the form:

$$\overbrace{C_1^1 C_2^1 C_3^1 \wedge C_1^2 C_2^2 C_3^2 \wedge \dots \wedge C_1^k C_2^k C_3^k}^{\text{conjunction of } k \text{ literal propositions}} \Rightarrow \overbrace{C_1^0 C_2^0 C_3^0}^{\text{conclusion}} \quad (7)$$

each literal made of m atomic concepts, $m = 3$ here

where a **concept** can be roughly understood as a **word vector** as in Word2Vec [20]. Each $C \in \mathbb{R}^d$ where d is the dimension needed to represent a single word vector or concept.

We use a “free” neural network (*ie*, standard feed-forward NN) to approximate the set of *all* rules. The **input** of the NN would be the state vector \mathbf{x} :

$$C_1^1 C_2^1 C_3^1 \wedge C_1^2 C_2^2 C_3^2 \wedge \dots \wedge C_1^k C_2^k C_3^k. \quad (8)$$

We fix the number of conjunctions to be k , with the assumption that conjunctions of length $< k$ could be filled with “dummy” (always-true) propositions.

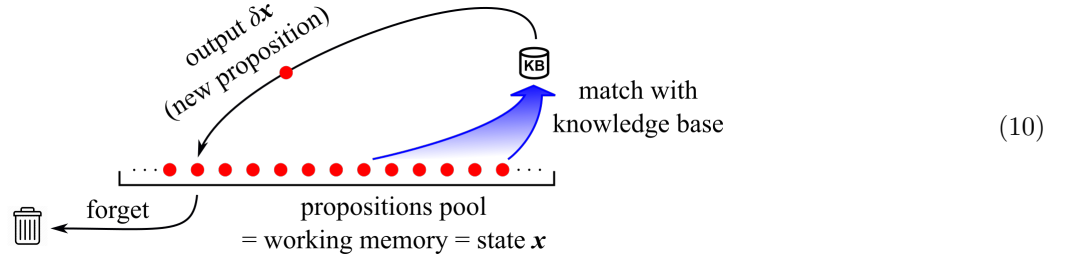
The **output** of the NN would be the conditional **probability** of an action:

$$P(\text{action} \mid \text{state}) := \pi(C_1 C_2 C_3 \mid \mathbf{x}). \quad (9)$$

Note that we don’t just want the action itself, we need the **probabilities** of firing these actions. The **Bellman update** of reinforcement learning should update the conditional probabilities over such actions (§??).

2.4 Structure of a logic-based AI system

Besides the intrinsic structure of a logic, the AI system has a structure in the sense that it must perform the following operations iteratively, in an endless loop:



- **Matching** — the $\boxed{\text{KB}}$ of rules is matched against the current state x , resulting in a (stochastically selected, *eg* based on ϵ -greedy) rule:

$$\boxed{\text{Match}} \quad (x \stackrel{?}{=} \boxed{\text{KB}}) : \mathbb{X} \rightarrow (\mathbb{X} \rightarrow \mathbb{P})$$

$$x \mapsto r$$
(11)

— In categorical logic, matching is seen as finding the **co-equalizer** of 2 terms which returns a **substitution** [4] [8] [9] [7]. The substitution is implicit in our formulation and would be *absorbed* into the neural network in our architecture.

— Matching should be performed over the entire **working memory** = the state x which contains k literals. This is combinatorially time-consuming. The celebrated **Rete** algorithm [19] turns the set of rules into a tree-like structure which is efficient for solving (11).

- **Rule application** — the rule is applied to the current state x to produce a new literal proposition δx :

$$\boxed{\text{Apply}} \quad r : \mathbb{X} \rightarrow \mathbb{P}$$

$$x \mapsto r(x) = \delta x$$
(12)

- **State update** — the state x is *destructively* updated where one literal $P_j \in x$ at the j -th position is **forgotten** (based on some measure of attention / interestingness) and over-written with δx :

$$\boxed{\text{Update}} \quad x = (P_1, P_2, \dots, P_j, \dots, P_k) \mapsto (P_1, P_2, \dots, \delta x, \dots, P_k)$$
(13)

All these operations are represented by functions parametrized by some variables Θ and they must be made *differentiable* for gradient descent.

2.5 Commutativity of logic conjunctions

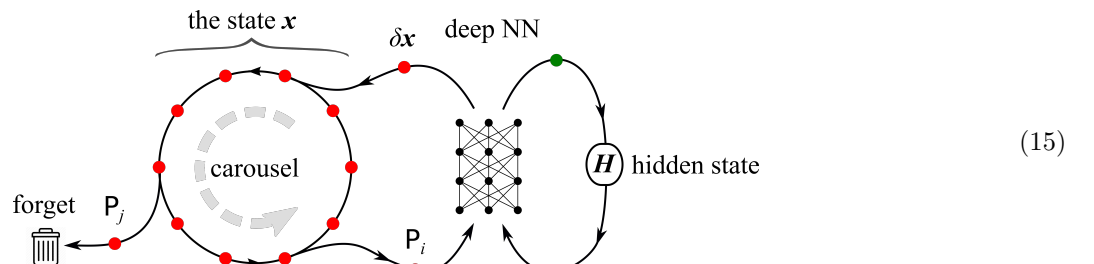
One basic characteristic of (classical) logic is that the conjunction \wedge is **commutative**:

$$P \wedge Q \Leftrightarrow Q \wedge P.$$
(14)

This restriction may significantly reduce the size of the search space. If we use a neural network to model the deduction operator $\vdash : \mathbb{P}^k \rightarrow \mathbb{P}$, where \mathbb{P} is the space of literal propositions, then this function should be **symmetric** in its input arguments.

I have considered a few solutions to this problem, including an algebraic trick to build “symmetric” neural networks (but it suffers from combinatorial inefficiency), and using Fourier transform to get a “spectral” representation of the state, which remained rather vague and did not materialize.

As of this writing ¹ I have settled on the “carousel” solution: All the propositions in working memory will enter into a loop, and the reasoning operator acts on a hidden state $H = \bullet$ and one proposition $P_i = \bullet$ at a time:



¹ Convolutional NNs are only *translation*-invariant. The Transformer [16] architecture is *equivariant* under permutations (meaning permuted inputs give permuted outputs), but it implicitly contains a recurrence similar to ours.

Notice that the working memory \mathbf{x} is itself a hidden state, so \mathbf{H} can be regarded as a *second-order* hidden state.

This architecture has the advantage of being simple and may be biologically plausible (the human brain’s working memory).

I believe in the maxim: *Whatever can be done in time can be done in space*. The diagram (15), when unfolded in time, can be expressed in this functional form:

$$\mathbf{F}_{\text{sym}}(\mathbf{P}_0, \dots, \mathbf{P}_k) = \mathbf{f}(\mathbf{P}_k, \mathbf{f}(\mathbf{P}_{k-1}, \mathbf{f}(\dots, \mathbf{f}(\mathbf{P}_0, \vec{\emptyset}))))). \quad (16)$$

\mathbf{F}_{sym} means that the function is invariant under the action of the symmetric group \mathfrak{S}_k over propositions. Such symmetric NNs have been proposed in [3] [2] [12] [13] [10] [11] [21].

2.6 Logic actions

The output of a logic rule is a proposition $\delta\mathbf{x} \in \mathbb{P}$, the space of propositions. This is a continuous space (due to the use of Word2Vec embeddings).

From the perspective of reinforcement learning, we are performing an action \mathbf{a} (the logic rule) from the current state \mathbf{x} . The neural network in (15) is a parametrized function \mathbf{F}_{Θ} that accepts \mathbf{x} and outputs $\delta\mathbf{x}$. We want to **gradient-descent** on Θ to get the optimal set of actions, ie, a **policy**.

To solve the RL problem, 2 main options are: **value-iteration** (eg Q-learning) and **policy-iteration**.

In **Q-learning**, we try to learn the action-value function $Q(\mathbf{a}|\mathbf{x}) \rightarrow \mathbb{R}$. The policy is obtained by choosing $\arg \max_{\mathbf{a}} Q(\mathbf{a}|\mathbf{x})$ at each step. In our logic setting, the action space $\mathbb{A} \ni \mathbf{a}$ is continuous. As is well known, if we use an NN to represent $Q(\mathbf{a}|\mathbf{x})$, the evaluation of $\arg \max_{\mathbf{a}}$ would be rather awkward, which is why Q-learning is widely seen as ill-suited for continuous actions.

It is much easier for the NN to directly output (a fixed number of) stochastic actions, thus avoiding the *curse of dimensionality*. Such a function is a **stochastic policy** $\pi(\mathbf{a}|\mathbf{x})$.

In other words, we can use a fixed number of Gaussian kernels (radial basis functions) to approximate the conditional probability distribution of π over \mathbb{A} . For each state \mathbf{x} , our NN outputs a probabilistic *choice* of c actions. So we only need to maintain c “peaks” given by Gaussian kernels. Each peak is determined by its mean $\delta\mathbf{x}_i$ and variance σ . Both parameters are to be learned.

The size of the FFNN in (15) seems well within the capacity of current hardware.

2.7 Forgetting uninteresting propositions

In (10) and (15), some propositions need to be forgotten based on some measure of **interestingness**. One way to measure interestingness is through the value function of a state, $V(\mathbf{x})$, where \mathbf{x} consists of propositions \mathbf{p}_i . Suppose that $V(\mathbf{x})$ is learned by a neural network, then it may be possible to extract (backwards) the weight by which a proposition \mathbf{p}_i contributed to the value V . For this to work, the function $V(\mathbf{x})$ should be deliberately **regularized** so that it would *generalize broadly*. Notice that $V(\mathbf{x})$ would also be *symmetric* in the \mathbf{p}_i ’s so it would have the architecture of (15). $V(\mathbf{x})$ is very similar to $Q(\mathbf{a}|\mathbf{x})$ but should be learned separately because they serve different purposes.

3 Remaining work

- Replace the recurrent NN architecture with symmetric NNs
- In this minimal architecture there is no **episodic memory** or **meta-reasoning** ability, but these can be added to the architecture and are not bottleneck problems. For example, meta-reasoning can be added via turning the input to introspection.
- Implementation of the system is currently under way.

Acknowledgements

I would like to thank the following people for invaluable discussions over many years: Ben Goertzel, Pei Wang (王培), Abram Demski, Russell Wallace, SeH, Jonathan Yan, and others. Also thanks to all the university professors and researchers in Hong Kong (especially in the math departments, and their guests), strangers who taught me things on Zhihu.com (知乎), Quora.com, and StackOverflow.

References

1. Andreka, Nemeti, and Sain. *Handbook of philosophical logic*, chapter Algebraic logic, pages 133–247. Springer, 2001.
2. de Bie, Peyré, and Cuturi. Stochastic deep networks. 2019. <https://arxiv.org/pdf/1811.07429.pdf>.
3. Gens and Domingos. Deep symmetry networks. *Advances in neural information processing systems*, (27):2537–2545, 2014.
4. Goguen. What is unification - a categorical view of substitution, equation and solution. *Resolution of equations in algebraic structures*, 1: algebraic techniques:217–261, 1989.
5. Halmos. *Algebraic logic*. Chelsea, 1962.
6. Halmos. *Logic as algebra*. Math Asso of America, 1998.
7. Bart Jacobs. *Categorical logic and type theory*. Elsevier, 1999.
8. Lawvere. *Functorial semantics of algebraic theories*. PhD thesis, Columbia university, 1963.
9. Lawvere and Rosebrugh. *Sets for mathematics*. Cambridge, 2003.
10. Qi, Su, Mo, and Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. 2016. <https://arxiv.org/abs/1612.00593>.
11. Qi, Yi, Su, and Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. *Advances in Neural Information Processing Systems*, pages 5105–5114, 2017.
12. Ravanbakhsh, Schneider, and Poczos. Deep learning with sets and point clouds. 2016. <https://arxiv.org/abs/1611.04500>.
13. Ravanbakhsh, Schneider, and Poczos. Equivariance through parameter-sharing. 2017. <https://arxiv.org/abs/1702.08389>.
14. Tarski, Henkin, and Monk. *Cylindric algebras, Part I*. 1971.
15. Tarski, Henkin, and Monk. *Cylindric algebras, Part II*. 1985.
16. Vaswani, Shazeer, Parmar, Uszkoreit, Jones, Gomez, Kaiser, and Polosukhin. Attention is all you need. 2017. <https://arxiv.org/abs/1706.03762>.
17. Wikipedia. AIXI. <https://en.wikipedia.org/wiki/AIXI>.
18. Wikipedia. No free lunch theorem. https://en.wikipedia.org/wiki/No_free_lunch_theorem.
19. Wikipedia. Rete algorithm. https://en.wikipedia.org/wiki/Rete_algorithm.
20. Wikipedia. Word2vec. <https://en.wikipedia.org/wiki/Word2vec>.
21. Zaheer, Kottur, Ravanbakhsh, Poczos, Salakhutdinov, and Smola. Deep sets. *Advances in Neural Information Processing Systems*, (30):3391–3401, 2017.