



Question answering by humans and machines: A complexity-theoretic view

Jan van Leeuwen^{a,*}, Jiří Wiedermann^{b,1}

^a Department of Information and Computing Sciences, Utrecht University, Princetonplein 5, 3584 CC Utrecht, the Netherlands

^b Institute of Computer Science, Czech Academy of Sciences, Pod Vodárenskou věží 2, 182 07 Prague 8, Czech Republic

ARTICLE INFO

Article history:

Received 12 June 2018

Accepted 2 August 2018

Available online 16 August 2018

Communicated by Josep Díaz

Keywords:

Question answering

Computational complexity

Human agents

Cognitive automata

Background intelligence

QA-machines

Advice

Learning space

Pippenger's theorem

Turing machines

ABSTRACT

Question-answering systems like Watson beat humans when it comes to processing speed and memory. But what happens if we compensate for this? What are the fundamental differences in power between human and artificial agents in question answering? We explore this issue by defining new computational models for both agents and comparing their computational efficiency in interactive sessions.

Concretely, human agents are modeled by means of *cognitive automata*, augmented with a form of *background intelligence* which gives the automata the possibility to query a given Turing machine and use the answers from one interaction to the next. On the other hand, artificial question-answering agents are modeled by *QA-machines*, which are Turing machines that can access a predefined, potentially infinite knowledge base ('advice') and have a bounded amount of *learning space* at their disposal.

We show that cognitive automata and QA-machines have exactly the same potential in realizing question-answering sessions, provided the resource bounds in one model are sufficient to match the abilities of the other. In particular, polynomially bounded cognitive automata with background intelligence (i.e. human agents) prove to be equivalent to polynomially bounded QA-machines with logarithmic learning space. It generalizes Pippenger's theorem on the computational power of switching circuits (without background intelligence) to a foundational result for question answering in cognitive science. The framework reveals why QA-machines have a fundamental advantage.

© 2018 Elsevier B.V. All rights reserved.

[These] phenomena have to be observed before being understood and explained with theories.

[Maurice Nivat [16]]

1. Introduction

In 2011, two experienced human players were beaten by IBM's supercomputer *Watson* at the game of *Jeopardy!* [2,4]. Can the question-answering capabilities of the human players and those of question-answering machines like *Watson* be placed in a rational perspective? In particular, can their potential be captured and compared in computational terms? In this paper we attempt to do so.

* Corresponding author.

E-mail addresses: J.vanLeeuwen1@uu.nl (J. van Leeuwen), Jiri.Wiedermann@cs.cas.cz (J. Wiedermann).

¹ The research of this author was partially supported by ICS CAS fund RVO 67985807 and CAS Programme Strategy 21.

In our approach to the problem, we make use of theoretical models for the human agents and Watson respectively, extending classical automata-theoretic conceptions for both of them. We then compare the relative capabilities of human agents and Watson-type machines in question-answering sessions (QA-sessions). In these sessions, sequences of questions are asked that are coded as words of length (at most) n , for some fixed but arbitrary $n \in \mathbb{N}$. We assume that all questions (input words) and answers (output words) are written over a fixed alphabet Σ with $|\Sigma| = \sigma \geq 2$.

In order to compare the qualities of humans and machines in question answering, we need suitable models that express their different ways of dealing with input, output, memory, discovery, and reasoning. The options for modeling the underlying processes in a fitting architecture for theoretical analysis are plenty. We will use a combination of abstract computational models that all have a long standing in the philosophy of mind and in cognitive science [15,18], with extensions as studied in modern complexity theory [1].

A guiding principle in our modeling effort is the *duality* between human players and question-answering machines, in the following sense. One may view human players as agents that consist of a knowledge-laden brain and a 'background intelligence' that can answer queries on call. On the other hand, one may view question-answering machines as agents that consist of an artificial 'foreground intelligence' to answer queries and a knowledge base in the background that contains all knowledge that it can use.

Finally, we embed some form of 'learning' in the models. Concretely, we will assume that a suitable facility is in place for preserving some amount of data in passing from one question to a next within a session. The human agent could preserve his 'state of mind' after a question, and a question-answering machine could preserve the contents of a kind of 'persistent memory'. No data is preserved in either model across session boundaries. The question arises how these facilities enrich the question-answering capabilities of each of the models, and how they compare.

1.1. Models: overview and intuitions

We first give an overview of the models as we envision them.

Modeling the human agent Cognitive science often models the information processing abilities of humans by means of finite-state automata with some kind of interactive quality. As processing qualities are expected to evolve as inputs get longer, we will assume that a separate automaton C_n is specified and called into action for each input length $n \geq 1$. 'Intelligence' is provided to each automaton, by allowing each C_n to access the computational support of a shared Turing machine M that acts as subliminal source of knowledge (or oracle) to all of them.

With these conceptions, a human agent can now be modeled as a *cognitive array* (CA) consisting of the following parts: a certain family $C = \{C(n)\}_{n \geq 1}$ of 2-way finite-state automata, and a single background Turing machine M_C that can be called by the automata. We assume that each CA C comes with a function s such that for each n , $s(n)$ bounds the size of its component $C(n)$.

Given a cognitive array C , a QA-session of C on inputs ('questions') of length n proceeds by *iterating* $C(n)$ on the successive inputs of the session. While $C(n)$ is processing an input question it can issue queries to M_C , using an oracle tape that is 'left-to-right write-only' during the process. The oracle tape is erased before a next question in the session is taken up, but $C(n)$ is assumed to be *persistent* in the sense that it continues in the state it was in after answering its most recent question.

Modeling question-answering machines In computing, the information processing abilities of question-answering machines such as Watson tend to be modeled by means of augmented Turing machines that have some kind of (potentially infinite) 'knowledge base' at their disposal. As the complexity of questions evolves as inputs get longer, we assume that for each $n \geq 1$ an appropriate string of data $B_n \in \Sigma^*$ is specified which is called into memory from the knowledge base whenever an input of length n is entered. Further data may accumulate during a session in a separate 'learning space' and may be passed on as newly discovered knowledge between successive queries. When a session ends, the learning space is cleared.

With these broad ideas we can model an artificial agent like Watson as a *question-answering machine* (QA-machine) consisting of the following parts: a multi-tape Turing machine W , a tape that serves as a learning space, and a knowledge base $B = \{B(n)\}_{n \geq 0}$. During a QA-session, W iterates on a series of questions of some fixed length n . Every time a question has been answered and a next one is presented in the session, W re-initializes but it does so while preserving some information that was gained in the handling of the previous question (in its learning space). This gives the model the capability to 'learn' during a session.

To measure their effect, we assume that every QA-machine W comes with functions b and l such that $|B(n)| \leq b(n)$ for all n and learning space in QA-sessions on inputs of length n remains bounded by $l(n)$.

Modeling Question Answering We aim to compare human agents and their Watson-type competitors, by focusing on the essence of their question-answering abilities. We will not consider the mismatch in speed nor any competitive interaction that would allow one party to step in if another one is slow or fails an answer. We will not even check whether questions adhere to a specific domain description or whether the answers that are given are 'correct', although this could well be added to our model. We focus entirely on *comparing* the fundamental question-answering abilities of the different agents. Can one agent (e.g. the human player) *in principle* come up with the same answers as the other agent (Watson)? Can one agent simulate the other and if so, what are the scaling effects involved? To make this precise, several more details have to be agreed upon. For a background on the theories of question answering, see e.g. [13,14].

Definition 1. A QA-session on inputs of length n consists of a finite sequence of inputs u_1, \dots with $u_i \in \Sigma^n$ ($i \geq 1$) and corresponding outputs v_1, \dots with $v_i \in \Sigma^*$ such that for each i , v_i is output by the machine on input u_i after it has completed processing u_1, \dots, u_{i-1} . A QA-session is $h(n)$ -bounded, for some non-decreasing function h , if the length of all outputs in the session is bounded by $h(n)$.

Definition 2. Given a CA or QA-machine, the process of answering any individual question (from input to corresponding output) in a session is called a *run* of the respective automaton or machine in the session.

In order to make sure that CA's and QA-machines always answer, we assume for both that the runs on individual questions of length n in a session are implicitly or explicitly bounded by some time-bound $f(n)$. We assume that $f(n) \geq n$, as the entire input must at least be read. If $f(n)$ is bounded by a polynomial, we say that the corresponding CA or QA-machine is 'polynomial-time bounded'. We say that a CA or QA-machine is *polynomially bounded* if all its respective parameters are polynomially bounded.

1.2. Discussion

The models we use are potentially very powerful, computationally. For example, having a separate finite-state automaton $C(n)$ for each input length n , already gives CA's the processing power of general 'resource-bounded' families of circuits [1]. Also, QA-machines build on Turing machines with bounded *advice* [10], the latter being an oracle facility similar to the notion of knowledge base as we defined it. CA's and QA-machines have several more features that we need for our modeling purposes.

Circuits and Turing machines with advice are common to non-uniform complexity theory, and their *duality* is well-known in this theory [1]. For example, a result known as Pippenger's Theorem asserts that a set is 'recognizable' by a polynomial-size bounded family of circuits if and only if the set can be recognized by a polynomial-time bounded Turing machine with polynomial-length bounded advice [1,10,17]. (For a more general formulation, see e.g. [20], Section 2.3.) This raises the interesting question whether a more general duality result may be achieved, notably for the models we have defined.

The iterative behavior of CA's and QA-machines during QA-sessions gives these models the ability to pass information from previous inputs to the next, either by state (in the case of CA's) or through learning space (in the case of QA-machines). This makes the behavior of the models dependent on everything that happened in the session before, which is not covered by the models in Pippenger's theorem as it stands.

Question-answering scenarios, in combination with the passing of information in between consecutive questions, have been considered before in studies of *persistent Turing machines* [5–8,11]. In these studies, persistence arises e.g. in modeling interactive computation or in computing on infinite streams of data. See also the studies of *interactive Turing machines with advice* [19]. We use question answering in a different way and in a different context, to compare the performance of different types of agents.

1.3. Outline

We compare the capabilities of human agents and Watson-type computers in QA-sessions by analyzing the models we designed for them. The models are described in detail in Section 2.

In Section 2 we first show that CA's and QA-machines become strictly more powerful when their 'capabilities for learning' grow. For CA's this corresponds to a growing 'brain size' $s(n)$, and for QA-machines to a growing size bound $l(n)$ on the available learning space.

In section 3 we develop the general theoretical results that show that CA's and QA-machines can simulate each other, within quantifiable resource bounds. This leads to the following theorem: *polynomially bounded cognitive arrays realize the same QA-sessions as polynomially bounded QA-machines with $O(\log n)$ -size bounded learning space*. The theorem and some further consequences are presented in Section 4.

The result establishes the desired duality result for cognitive automata with background intelligence. It also gives a powerful extension of Pippenger's theorem for plain (acyclic, memory-less) switching circuits to the complex context of cognitive automata with advanced features like persistent memory and background intelligence. In Section 5 we give some final conclusions.

2. Models in detail

In order to derive simulation results, we need a more complete description of the models. Recall that both cognitive arrays and QA-machines are designed to perform QA-sessions of successive questions (inputs) of length n , with n fixed throughout a session. In Section 2.1 we describe how cognitive arrays operate, in Section 2.2 we describe QA-machines in detail. Finally, in Section 2.3 we comment on the use of QA-sessions for characterizing the strength of question-answering agents in general.

2.1. Modeling human agents: cognitive arrays

The premise of cognitive arrays is that question answering is based on the knowledge that is stored in the ‘states’ of the human brain, with an inference mechanism provided by a Turing machine that acts as a *computable oracle*. Hence, the inference mechanism is potentially far more powerful than plain finite-state reasoning. The precise architecture of the model is as follows.

A cognitive array C consists of an infinite sequence of components $\{C(n)\}_{n \geq 1}$ and one background multi-tape Turing machine $M = M_C$. Each $C(n)$ is a finite-state machine, with $s(n)$ states and a special capability to create $q(n)$ -size queries on an internal oracle tape and pose them to M . Whenever M is queried, we assume that it always responds ‘super fast’ (i.e. in one step).

If a QA-session starts and its first input of length n appears, C initializes $C(n)$ and delegates the further handling of the session to this automaton. Suppose $C(n)$ has handled zero or more inputs of the session and now receives the next question, say x , on its input tape ($|x| = n$). $C(n)$ continues in the state it was in after completing the answering of the preceding input, but its oracle tape is emptied before it starts its run on x . Now $C(n)$ acts as follows.

First of all, $C(n)$ processes x like a *two-way* automaton, with state transitions determined by the current input or oracle signal and the current state. (In *one-way* CA’s, this is limited to a one-way behavior on x .) The oracle signal can be ‘yes’ or ‘no’ and is received from a query request to machine M . (How this works is explained below.) In each action, also some output may be written to its one-way oracle tape and to a one-way output as well.

During its run on an input x , $C(n)$ can query the fixed background Turing machine M . To this end, $C(n)$ has its oracle tape of at most $q(n)$ positions, which is write-only and *written from left-to-right*. (When the oracle tape is full, an end marker will prevent it from overflowing.) When $C(n)$ gets into an oracle state, M is queried with the contents of the oracle tape, and $C(n)$ moves according to the answer it gets (‘yes’ or ‘no’) to a new state. The oracle tape is *not* emptied but keeps being written to by adding things to the right. M does not retain any information from one query to a next.

This continues until the machine stops, completing its answer to x on the output. When $C(n)$ is given a new input x' with $|x'| = n$, it *continues in the state it was in* but with an *emptied* oracle tape (i.e. the tape contents does not transfer). This continues input after input, as long as the session lasts. We assume that every run of $C(n)$ on an input in the session n is $f(n)$ -time bounded.

When M is queried as an ‘oracle’ we assume that it gives its answer like an oracle ‘in one step’. However, for consistency we require that M signals ‘yes’ precisely when it halts on the query in an accepting state, and ‘no’ when it halts in a rejecting state. Let $r(n)$ be the maximum time that M computes in any halting computation on queries of size at most $q(n)$. (Clearly, $r(n)$ is non-computable in general.)

Definition 3. A cognitive array C is called a $\langle s(n), q(n), r(n), f(n) \rangle$ -array, if C consists of a sequence of 2-way finite-state components $\{C(n)\}_{n \geq 1}$ and a background Turing machine M which conform to the parameters specified above: $s(n)$ bounds the number of states of $C(n)$, $q(n)$ bounds the length of the queries $C(n)$ can ask to M , $r(n)$ is the maximum time M needs in any halting computation on queries posed by $C(n)$, and $f(n)$ bounds the time of any run by $C(n)$, all on inputs of length n .

Note that the representation of $C(n)$ with a full transition table requires $O(\log |\Sigma| \cdot s(n) \log s(n))$ bits.

2.2. Modeling Watson-type agents: QA-machines

QA-machines are based on the premise that question answering is driven by an inference mechanism modeled as a Turing machine, supported by the retrieval of information from a large external knowledge base and a learning facility. The architecture of QA-machines is as follows.

A QA-machine consists of a front-end multi-tape Turing machine W with the following features: (a) W can access a knowledge base B containing strings $B(n) \in \Sigma^*$ of size $b(n)$ (at most one for each $n \geq 1$), and (b) W has a learning space (i.e. tape) of size $l(n)$ on which it can accumulate some information for later use in a given session. In addition, W has its usual set of working tapes as a multi-tape machine.

In a QA-session, W operates on successive inputs x with $|x| = n$ in the expected way, each time waiting until the answer to a question has been completed on its one-way output before a next input is taken in. While it operates on an input of length n , W can access its knowledge base to retrieve the stored item $B(n)$ from it with $|B(n)| \leq b(n)$. This information acts as an *advice* [10] that stays with the machine for the rest of the session. The string $B(n)$ holds all the relevant information which W can extract from its knowledge base and use in a single session.

W is not just a machine with ‘advice’. W also has a special tape of size at most $l(n)$ that acts as a learning space. The learning space is a *persistent* memory that keeps its contents in between different runs of the machine on queries of length n . The space allows W to accumulate some additional knowledge during a QA-session, e.g. by storing useful facts that it learns from its inputs. All other working tapes are cleared when a new run begins.

A QA-machine basically gets all its knowledge from the knowledge base B , and from what it learns from its query processing and can keep in its learning space after a run. Clearly, the QA-session thus depends on the historical information W builds up in its learning space, hence on the particular chain of inputs it is provided during a session. A different chain

of inputs might give different answers to identical queries later on, as one expects from a *learning machine*, although the accumulated data is only used during a single session. We assume that runs of W always halt within a given time-bound $f(n)$, implying that questions are always answered in finite time.

Definition 4. A QA-machine is called a $\langle b(n), l(n), f(n) \rangle$ -machine, if it consists of a front-end Turing machine W and a knowledge base B with $B(n) \in \Sigma^*$ ($n \geq 1$) which conform to the parameters specified above: $b(n)$ bounds the length of $B(n)$, $l(n)$ bounds the size of the learning space throughout a session, and $f(n)$ bounds the time of any run by the machine, all on inputs of length n .

We usually ‘name’ a QA-machine by the name of its Turing machine. Note that outputs in a QA-session with a $\langle b(n), l(n), f(n) \rangle$ machine are necessarily $f(n)$ -bounded.

2.3. Question answering: on the power of learning space

Question answering is a standard type of information retrieval and filtering. We defined QA-sessions only abstractly. In a session, questions of a specified complexity (length) are to be supplied by an external party, where a next question is supplied only if the answer to the current question has been output. We do not concern ourselves with the particular way the external party determines its next question.

In our present models we have amplified the different ways answering agents may store and act with information. While QA-sessions were designed primarily as a means to compare the strength of different answering agents, one could as well use QA-sessions to characterize the potential of any one answering party in particular. Let \mathcal{A} be any answering agent.

Definition 5. The QA-script of \mathcal{A} is the set $R(\mathcal{A})$ of all QA-sessions that are realized using \mathcal{A} as the answering agent.

Given a family of answering agents, one may study the class of all QA-scripts that can be realized by these agents, in the same way as classes of formal languages are studied in relation to the automata that accept them [9]. Some aspects of this are studied in [5–7,10]. We do not pursue this question here but only highlight the effect that various *model parameters* can have on the classes of QA-scripts that can or cannot be realized.

The quality of the question answering during a session depends on two factors: one is the intrinsic computational ability of the question-answering agent by itself, another is the facility by which information can be passed on from the answering of one question to the next. This information can be viewed as a ‘residue’ from preceding answers that can help in future questions in the same session [12].

Focusing on the latter, for example, it is intuitive that the more information can be passed, the more ‘powerful’ the question answering during a session becomes. We give a concrete example for QA-machines. The argument exploits a technique from automata theory involving crossing sequences [9].

Theorem 1. For any computable function $L : \mathbb{N} \rightarrow \mathbb{N}$, there exists a QA-script \mathcal{S} that can be realized by a QA-machine within a learning space of size $L(n)$, but not by any QA-machine with a learning space of size $L'(n)$ with $\inf_{n \rightarrow \infty} \frac{L'(n)}{L(n)} = 0$.

Proof. Let $n \in \mathbb{N}$ be arbitrary, and set $k = \lceil \frac{L(n)}{n} \rceil$. (The value of k can be computed by a QA-machine anytime it needs it.) We will consider sessions with questions of length n . Define a question-block to be any sequence of k questions (of length n each) of the form $uv_0, v_1, \dots, v_{k-1}$ where u is a fixed default string of length $kn - L(n)$. Obviously $|v_0| + \dots + |v_{k-1}| = L(n)$. Given a question-block $\alpha = uv_0, v_1, \dots, v_{k-1}$ we say that $v_0v_1 \dots v_{k-1}$ is the string stored in it. This gives a unique correspondence between question-blocks and strings of length $L(n)$ over Σ .

View strings of length $L(n)$ over Σ as σ -ary numbers. Let $l = \sigma^{L(n)}$ and let $\alpha_0, \dots, \alpha_{l-1}$ be the question-blocks corresponding to the complete list of length- $L(n)$ strings over Σ in enumeration order. We use these blocks to define the QA-script \mathcal{S} consisting of sessions which adhere to the following rules. Our aim is to view any session as a sequence of blocks of k questions each. First assume that the number of questions in a session is a multiple of k .

- If the first block of k questions equals α_0 , then the corresponding block of outputs consists of $k - 1$ times *no* followed by one *yes*, otherwise it consists of k times *no*.
- If the last block of k answers consisted of k times *no*, then every next block of k questions is answered by k times *no*.
- Suppose the last block of k answers consisted of $k - 1$ times *no* followed by one *yes* and that this last block of k questions equaled α_i (some $0 \leq i \leq l - 1$). Let $i' = (i + 1) \bmod l$. If the next block of k questions corresponds to $\alpha_{i'}$, then the corresponding block of outputs consists of $k - 1$ times *no* followed by one *yes*, otherwise it consists of k times *no*.

If the number of questions in a session is not a multiple of k , then apply the rules to the largest prefix with a multiple of k questions and answer *no* to the remaining (fewer than k) questions.

Claim. S is realized by a QA-machine with a learning space of size $L(n)$.

Proof. A QA-machine W can be designed which answers in all sessions exactly according to the rules. The idea is to work through blocks of k consecutive questions, and pass a copy of the desired α_j from question to question in learning space. By checking that the k consecutive parts of α_j correspond exactly (from left to right) to the k consecutive questions in the input, W can verify that the desired match occurs. The process is initialized by starting with the empty (default) learning space. If a match is made, the contents of learning space is incremented by one as a number (modulo l), and W is ready for the next block of k inputs. If no match is made, then the *no* answer perpetuates from then onward. Note that for passing a question-block in learning space, it suffices to pass only the ‘significant’ part of the block that corresponds to the size- $L(n)$ string stored in it. By marking segments, one can keep track of which part of the block must match the current input question. This can all be encoded within $L(n)$ workspace. \square

Let $L'(n) : \mathbb{N} \rightarrow \mathbb{N}$ be an arbitrary function, not necessarily computable, with $\inf_{n \rightarrow \infty} \frac{L'(n)}{L(n)} = 0$.

Claim. S cannot be realized by any QA-machine with a learning space of size $L'(n)$.

Proof. Suppose there was a QA-machine W' that realized S with a learning space of size $L'(n)$. For some n to be chosen later, consider how W' would act in a session with questions $\alpha_0, \alpha_1, \dots, \alpha_{l-1}, \alpha_0$. Clearly, W' must answer with repeated blocks of $k-1$ times *no* followed by one *yes*.

Let W' have a working alphabet on its learning tape of δ symbols. Then there can be at most $\delta^{L'(n)}$ different learning tape contents when W' passes from some block α_i to the next ($0 \leq i \leq l-1$). By assumption on L' there must be infinitely many values of n such that $L'(n)/L(n) < \log \sigma / \log \delta$. For any such value of n we have $\delta^{L'(n)} < l$.

Let n be any value with this property. By the pigeonhole principle, there are indices i, j with $0 \leq i < j \leq l-1$ such that W' passes from α_i to α_{i+1} and from α_j to $\alpha_{j'}$ with the same learning tape contents, where $j' = (j+1) \bmod l$. But, then we can omit the questions from α_{i+1} to α_j and W' must answer consistently with repeated blocks of $k-1$ times *no* followed by one *yes* even though the series of questions has been altered to $\alpha_0, \dots, \alpha_i, \alpha_{j'}, \dots$ which clearly should have led to k times *no* when $\alpha_{j'}$ is processed. Contradiction. \square

This completes the proof. \square

A similar result can be shown for cognitive arrays, using the size bound $s(n)$ of the state sets as an indicator: as $s(n)$ grows, the more powerful the question answering by the cognitive array can be. This is consistent with the general idea that *the bigger an agent's brain (or memory), the smarter it can be*.

3. Cognitive arrays versus question-answering machines

The key issue of this paper can be now phrased as follows: *are human agents and Watson-type QA-machines able to realize the same QA-scripts?* And, if the answer is positive, are there any appreciable differences in effectiveness between the two, not counting the obvious difference in speed?

We will show that the question-answering power of cognitive arrays corresponds exactly to that of suitably configured QA-machines, and vice versa. We prove this by showing that one model can realize the same QA-sessions as the other, with the right parameter settings. We assume that all models work with a fixed alphabet Σ with $|\Sigma| \geq 2$. Constant factors will be suppressed in the various bounds.

3.1. Simulating human agents by QA-machines

We first show how QA-machines can simulate the capabilities of cognitive arrays of any given setting.

Theorem 2. Let Φ be any QA-session. If Φ is realized by a $\langle s(n), q(n), r(n), f(n) \rangle$ -array, then the session can also be realized by a $\langle s(n) \log s(n) + q(n), \log s(n), q(n) + f(n)(r(n) + s(n) \log s(n)) \rangle$ -machine.

Proof. Let C be an $\langle s(n), q(n), r(n), f(n) \rangle$ -array and let $M = M_C$ be the multi-tape Turing machine it uses as an oracle. Consider the QA-session Φ on inputs of length n . We design a QA-machine W that simulates $C(n)$ as follows.

First of all we need to deal with the simulation of M on queries of length up to $q(n)$. Define:

$\gamma(n) =$ “the lexicographically least string of length $\leq q(n)$ that leads M to perform its longest halting computation over all strings of length $\leq q(n)$.”

We assume $\gamma(n)$ is always defined, but give it a default value in those cases in which it isn't. If W knows $\gamma(n)$, then it can simulate M in finite time, using the computation on $\gamma(n)$ as a yardstick to decide when to terminate any computation of M on a query of length $q(n)$ or less. It is immaterial for our purposes that $\gamma(n)$ may be non-computable.

We now define W in a number of steps. To begin with, we give it a knowledge base B with

$$B(n) = \text{“the description of } C(n) \text{ followed by } \gamma(n)\text{”}.$$

Given any input x with $|x| = n$, W first retrieves $B(n)$ and unpacks the description of $C(n)$ and $\gamma(n)$ (encoded e.g. in binary) in $s(n) \log s(n) + q(n)$ time and then starts to simulate the execution of $C(n)$ on x . Whenever an oracle call to M is made, W simulates the call in finite time with the computation on $\gamma(n)$ to bound it. When the simulation of $C(n)$ ends, W stores an encoding of the ending state of $C(n)$ in its learning space. This enables W to start the simulation of $C(n)$ on a next input in the session correctly, namely in the state which it has recorded in its learning space. Repeating this, W precisely realizes Φ .

We conclude by noting that $|B(n)| = b(n) \leq s(n) \log s(n) + q(n)$ and $l(n) \leq \log s(n)$. The simulation of $C(n)$ takes $O(s(n) \log s(n))$ per step, to retrieve the right move from the encoded transition table, plus an additional $O(r(n))$ time if the oracle is called. \square

Note in the given proof that we do not really need to ‘know’ $r(n)$ or $f(n)$ in the construction as we merely used them for bounding the time-complexity of W .

Corollary 1. *If a QA-session Φ is realized by a $\langle s(n), q(n), r(n), f(n) \rangle$ -array with $r(n) \geq s(n) \log s(n)$, then the session can also be realized by a $\langle s(n) \log s(n) + f(n), \log s(n), f(n)r(n) \rangle$ -machine.*

Proof. By the operation of $C(n)$ we may assume that $q(n) \leq f(n)$. Assuming that $r(n) \geq s(n) \log s(n)$ we obtain the desired result, suppressing constant factors. \square

Corollary 2. *If a QA-session Φ is realized by a $\langle s(n), q(n), r(n), f(n) \rangle$ -array with $r(n) \geq s(n) \log s(n)$, then Φ can also be realized by a $\langle s(n)(nq(n) + \log s(n)), \log s(n), ns(n)q(n)r(n) \rangle$ -machine.*

Proof. By the operation of $C(n)$ we may assume that $f(n) \leq n \cdot s(n) \cdot q(n)$. After all, in between any two subsequent writes to its oracle tape, the finite-state machine $C(n)$ cannot make more than $n \cdot s(n)$ moves, otherwise it would repeat an input position/state pair and get into a loop. Now apply Corollary 1 to get the desired bound. \square

3.2. Simulating Watson-type QA-machines by cognitive arrays

Next we reduce the question-answering intelligence of QA-machines to that of, suitably dimensioned, cognitive arrays.

Theorem 3. *Let Φ be a QA-session. If Φ is realized by a $\langle b(n), l(n), f(n) \rangle$ -machine, then the session can also be realized by a $\langle b(n) + 2^{O(l(n))}, b(n) + l(n) + f(n), f(n), b(n) + 2^{O(l(n))} + f(n) \rangle$ -array.*

Proof. Let W be a $\langle b(n), l(n), f(n) \rangle$ -machine, and let B be the knowledge base it uses. We show that a cognitive array C can be designed such that for every n , the automaton $C(n)$ encodes W ’s knowledge insofar as it is needed in processing inputs of length n and can simulate W during the session on successive inputs x as processed in Φ , using a suitable multi-tape Turing machine $M = M_C$ as background intelligence. We develop the description of $C(n)$ and the states it needs in stages, for any n .

Phase I: Preparations.

First we need a representation of the initial contents of the learning tape. Let A be the complete σ -ary tree of depth $l(n)$, with the edges at every node labeled by the different symbols of Σ . Each node of A thus represents a possible content of W ’s learning space, by collating the edge labels on the path from the node to the root. With this correspondence we identify each internal node of T with a possible tape contents of the learning space, a so-called A -state. This leads to some $\sigma^{l(n)}$ A -states.

Assume that $C(n)$ starts a new run in state α , corresponding to the contents of W ’s learning space at the end of the previous (simulation of a) run. $C(n)$ first reads and copies the new input x to its oracle tape, symbol by symbol. Next, it switches over to a different mode and moves from state α up the tree towards the root, copying the symbols on the way to the oracle tape as well. The oracle tape now contains $\langle x, \vec{\alpha}, \dots \rangle$, where $\vec{\alpha}$ is the string corresponding to state α . Note that the total number of states to execute this phase so far is bounded by $O(\sigma^{l(n)})$.

Next, $C(n)$ needs to ‘internalize’ $B(n)$. To this end, we hard-wire $B(n)$ into $C(n)$, by means of a row of (up to) $b(n)$ states. After completing the previous steps, $C(n)$ starts at the beginning of the row of states and moves from left to right, outputting the corresponding symbols to its oracle tape again. When the last state in the row is passed, $C(n)$ moves to the next phase.

Phase II: Generating output

The oracle tape now contains $\langle x, \vec{\alpha}, B(n), \dots \rangle$. In order to generate output, $C(n)$ moves to yet another set of special states, the Y -states. In its Y -states, $C(n)$ repeatedly adds a 1 to its oracle tape and then queries its oracle $\sigma + 1$ times

as follows. When the oracle tape contains $\langle x, \vec{\alpha}, B(n), 1^y \rangle$, the following oracle queries are asked, by moving through subsequent states:

- for each $\sigma \in \Sigma$: “does W print out a symbol σ in the y -th step of its computation on input x , with $\vec{\alpha}$ as its initial learning tape contents and $B(n)$ as the data in its knowledge base for this session?”
- and finally: “does W terminate in the y -th step of its computation on input x , with $\vec{\alpha}$ as its initial learning tape contents and $B(n)$ as the data in its knowledge base for this session?”

Clearly we need $\sigma + 1$ Y -states to perform the queries. Whenever the answer to a σ -query is ‘yes’, $C(n)$ outputs σ and continues its cycle. If the answer to the termination-query is ‘no’, $C(n)$ writes another 1 to its oracle tape and continues with a new cycle. If the answer is ‘yes’, $C(n)$ ends the cycling and moves to the next phase. At this point $C(n)$ has generated the output W would have generated. Note that the oracle machine must be able to answer its queries in the time W needs for its computations, i.e. within $f(n)$ steps.

Phase III: Wrapping up.

Note that the oracle tape now holds $\langle x, \vec{\alpha}, B(n), 1^K \rangle$, where K ($K \leq f(n)$) is the number of steps W takes to complete its output on x . To wrap up, $C(n)$ must reconstruct the contents of the final learning tape and move to the corresponding A -state, to prepare for the next input. To this end, $C(n)$ moves to its Z -states. The Z -states are a replica of the A -states, with transitions between them now corresponding to a top-down, level-wise enumeration order of the nodes of tree A . Starting at the root and moving through the Z -states one by one in enumeration order, level after level, $C(n)$ adds each time yet another 1 to the oracle tape and queries the oracle with $\langle x, \vec{\alpha}, B(n), 1^K, 1^z \rangle$ on its tape with the following question:

- “does W terminate with a contents corresponding to the z -th node in the enumeration in its learning space, in the computation on input x , with $\vec{\alpha}$ as initial learning tape contents and $B(n)$ as the data in its knowledge base for this session?”

Once the oracle says ‘yes’, we have in fact reached the node in the A -tree that corresponds to the contents of the learning space that must carry over and $C(n)$ can move to the corresponding A -state α' , to prepare for a next run in the session. Note that we do not need K for these last queries, although we have its value on the oracle tape.

For the final accounting, observe the following. The required oracle is easily seen to be implementable as a uniform TM M for all n . Considering the queries we ask, M always halts and does so within time $f(n)$ by definition. A single run of $C(n)$ takes time in the order of $n + l(n) + b(n) + \sigma f(n) + \sigma^{l(n)} = O(b(n) + 2^{O(l(n))} + f(n))$, with oracle queries taking ‘one step’. We assume that $f(n) \geq n$. \square

Interestingly, the construction automatically yields a one-way cognitive automaton that simulates the QA-machine within the given complexity bounds.

Corollary 3. *Let Φ be a QA-session. If Φ is realized by a $\langle b(n), l(n), f(n) \rangle$ -machine, then it can also be realized by a one-way $\langle b(n) + 2^{O(l(n))}, b(n) + l(n) + f(n), f(n), b(n) + 2^{O(l(n))} + f(n) \rangle$ -array.*

Proof. Considering the construction above, we see that Phase I is the only phase in which the (next) input of a session is concretely read by $C(n)$. It does so in a straight one-way manner. \square

Note once again that we did not need $f(n)$ in the construction but only used it in bounding the time-complexity of $C(n)$. We may assume w.l.o.g. that $f(n)$ dominates both $b(n)$ and $l(n)$, as W must read both the query and the learning tape or can be made to do so. This leads to the following corollary.

Corollary 4. *If a QA-session Φ is realized by a $\langle b(n), l(n), f(n) \rangle$ -machine and $f(n) \geq \max\{b(n), l(n)\}$, then the session can also be realized by a (one-way) $\langle b(n) + 2^{O(l(n))}, f(n), f(n), f(n) + 2^{O(l(n))} \rangle$ -array.*

4. Some consequences

It is tempting to hypothesize that there are no principal differences in intellectual capability between a human agent and a question-answering machine like Watson, except for those caused by the sheer unlimited memory and processing speed of the latter. In the models we presented it becomes meaningful to address this question. Can the perceived advantage of question-answering machines as artificial brains be made tangible?

The simulation results proved in Section 3 are tedious but they show that in principle, cognitive arrays and QA-machines are *equally powerful*. The models can be matched to each other on any QA-script, provided they are tuned as in the simulations. It certainly isn't the whole story. For example, do the simulations point to any differences in effectiveness between (the models of) the human agents and Watson?

4.1. Overview

Our comparison between human agents and machines focused entirely on their capabilities in question answering, in relation to the ‘system parameters’ we distinguished for them. We abstracted away from physical details and allowed for perfect generality in the scenario of QA-sessions and scripts. This seems to be a new approach to the analysis of question answering in general.

The automata-theoretic models we used attempted to capture the functionalities of human agents and question-answering machines. The persistence of states and learning space, respectively, gives an interactive feature to the models which substantially extends the regime of traditional machine models, which at best have ‘advice’ only and no persistence ([3,10]). The results of Section 3 can be summarized as follows, where ‘simulation’ refers to the capability of realizing equal QA-sessions and the natural assumptions on $f(n)$ are omitted.

- From human agents to QA-machines: every cognitive $\langle s(n), q(n), r(n), f(n) \rangle$ -array with $r(n)$ satisfying $r(n) \geq s(n) \log s(n)$ can be simulated by a $\langle s(n) \log s(n) + f(n), \log s(n), f(n)r(n) \rangle$ -machine.
- From QA-machines to human agents: every $\langle b(n), l(n), f(n) \rangle$ -machine can be simulated by a (one-way) cognitive $\langle b(n) + 2^{O(l(n))}, f(n), f(n) + 2^{O(l(n))} \rangle$ -array.

4.2. Extending Pippenger’s theorem

The duality between CA’s and QA-machines reminds of that between combinational circuits and Turing machines with advice studied in non-uniform complexity theory. One of the key observations in the latter field is Pippenger’s theorem: *a language is accepted by a family of circuits of polynomially bounded size if and only if the language is accepted by a polynomial-time bounded Turing machine with polynomially bounded advice*. For a more detailed formulation and proof we refer to [1,10,17,20].

CA’s and QA-machines are based on machine models which extend these basic counterparts in substantial ways. Nevertheless, it is a natural question whether Pippenger’s theorem might not in essence be an instance of a more general result, for the more general models we defined. Let $T(n) \geq n$.

Definition 6. A $\langle s(n), q(n), r(n), f(n) \rangle$ -array is said to be $T(n)^{O(1)}$ -bounded if all its parameters are. A $\langle b(n), l(n), f(n) \rangle$ -machine is said $T(n)^{O(1)}$ -bounded if its parameters $b(n)$ and $f(n)$ are.

An immediate consequence of the simulation results from Section 3 is the following result, generalizing Pippenger’s theorem.

Theorem 4. $T(n)^{O(1)}$ -bounded QA-machines with $O(\log T(n))$ -size learning space are equivalent to (one-way) $T(n)^{O(1)}$ -bounded cognitive arrays.

Proof. The result follows when we choose $l(n) = \log T(n)$ and all other specified parameters $T(n)^{O(1)}$ -bounded. \square

Corollary 5. $2^{O(g(n))}$ -bounded QA-machines with $O(g(n))$ -size learning space are equivalent to (one-way) $2^{O(g(n))}$ -bounded cognitive arrays.

Corollary 6. Polynomially bounded QA-machines with $O(\log n)$ -size learning space are equivalent to (one-way) polynomially bounded cognitive arrays.

4.3. Effectiveness

From the viewpoint of computational cognitive science, more conclusions may be drawn. In the simulation of a QA-machine by a CA, the key observation is that $s(n) \approx b(n) + 2^{O(l(n))}$. The term $b(n)$ reflects the fact that the background knowledge of the QA-machine must be encoded in the brain automaton of the CA, as expected. However, one also sees that $s(n)$ has a term that is exponential in the size of the learning space of the QA-machine. Examples like the one in Theorem 1 show that this bound is best possible in general, in the worst case. The exponential term reflects the fact that the brain automaton of the CA must have the circuitry in place for any possible contents of the learning space that may arise. (If a better estimate for the latter can be given, a better bound may result.)

Thus, whereas the question-answering capabilities of the two models may be made equal by suitable parameter tunings, it cannot be avoided that the size bound of the automata in a simulating CA is exponential in the size of the learning space of the QA-machine, aside from the ‘linear’ term $b(n)$. This is the unavoidable consequence of having a model in which the brain of the CA must be hard-wired in advance for the knowledge that the given QA-machine has at its disposal during a session, either from its knowledge base or by accumulation in learning space. In contrast, the simulation of a CA by a QA-machine takes a setting that remains polynomial in all parameters of the CA.

The observations give an indication of the difference in effectiveness between the human agents and Watson, aside from the obvious differences in speed. The difference seems to be caused entirely by the advantages of a Watson-type computer in storing vast amounts of potentially useful information, possibly exceeding any feasible bound on the information that can be stored in the human brain, either in advance or by learning during a session.

5. Conclusion

In this paper we gave a theoretical model for comparing the qualities of human agents and machines at *question answering*. The motivating premise was the bold idea that the human agents and the supercomputer Watson may be equally good at the game of Jeopardy! in principle, if their capabilities match. Part of the effort in this paper went into identifying suitable models and parameters, for a meaningful comparison.

We modeled human agents as cognitive arrays and Watson-type computers by QA-machines, and we observed their performance during QA-sessions. Cognitive arrays model the ‘circuits of the mind’ of the human agents, with the option to learn and gain additional information from a Turing machine. QA-machines model the query-answering abilities of Turing machines with predefined advice and a persistent learning space. The models are very powerful, if we allow them to be as general as defined. In particular, without any further constraints cognitive arrays and QA-machines can encode and make use of arbitrary knowledge bases. This gives the models potentially ‘super-Turing power’ as computational devices, unless one insists on uniform definability of the infinite parts of the models. The present analysis has not made use of any such assumption.

A comparison of the models shows that the power of QA-machines as we have modeled them can be captured by cognitive arrays, but at the price of having state sets that are large enough so as to encode what a QA-machine stores in its knowledge base and can accumulate in its learning space during a session. This suggests the true advantage of QA-machines, having more effective means for storing and maintaining this information. The comparison also suggests that the model of cognitive arrays might have to be changed, to deal with temporary knowledge in a different way than through fixed circuitry. This seems to lead straightforwardly to the notion of neural nets which can modify their states, their transitions and even their topology [21]. For the time being, QA-machines have the edge.

Our main result shows that the models are ‘equally powerful’, provided their parameters are set appropriately. In particular, it follows that cognitive automata of polynomial size are equivalent to polynomially bounded QA-machines with logarithmic learning space. This result extends Pippenger’s theorem from non-uniform complexity theory, originally relating the computational power of circuits and Turing machines with advice only. The extension may be seen as a foundational result for the theory of question-answering machines.

References

- [1] J.L. Balcázar, J. Diaz, J. Gabarro, Structural Complexity I, 2nd ed., Springer-Verlag, Berlin, 1995.
- [2] B.B.C. News, IBM’s Watson supercomputer crowned Jeopardy king, <http://www.bbc.co.uk/news/technology-12491688>, February 17, 2011.
- [3] C. Damm, M. Holzer, Automata that take advice, in: J. Wiedermann, P. Hájek (Eds.), Mathematical Foundations of Computer Science 1995, Proc. 20th International Symposium (MFCS’95), in: Lecture Notes in Comput. Sci., vol. 969, Springer-Verlag, Berlin, 1995, pp. 149–158.
- [4] D. Ferrucci, et al., Building Watson: an overview of the DeepQA project, AI Mag. (Fall 2010) 59–79.
- [5] D. Goldin, P. Wegner, Persistence as a Form of Interaction, Tech. Rep. CS-98-07, Dept. of Computer Science, Brown University, Providence RI, 1998.
- [6] D.Q. Goldin, Persistent Turing machines as a model of interactive computation, in: K.-D. Schewe, B. Thalheim (Eds.), Foundations of Information and Knowledge Systems, in: Lecture Notes in Comput. Sci., vol. 1762, Springer-Verlag, 2000, pp. 116–135.
- [7] D.Q. Goldin, S.A. Smolka, P.C. Attie, E.L. Sonderegger, Turing machines, transition systems, and interaction, Inform. and Comput. 194 (2) (2004) 101–128.
- [8] H. Hempel, M. Kimmritz, Persistent computations of Turing machines, in: O.H. Ibarra, B. Ravikumar (Eds.), Implementation and Application of Automata, 13th Int. Conference, CIAA 2008, in: Lecture Notes in Comput. Sci., vol. 5148, Springer-Verlag, 2008, pp. 171–180.
- [9] J.E. Hopcroft, J.D. Ullman, Formal Languages and Their Relation to Automata, Addison-Wesley Publ. Comp., Reading MA, 1969.
- [10] R.M. Karp, R.J. Lipton, Turing machines that take advice, Enseign. Math. 28 (1982) 191–209.
- [11] S. Kosub, Persistent Computations, Techn. Rep. 217, Institut für Informatik, Julius-Maximilians-Universität Würzburg, Würzburg, 1998.
- [12] W.G. Lehnert, Paradigmatic issues in cognitive science, in: W. Kintsch, J.R. Miller, P.G. Polson (Eds.), Methods and Tactics in Cognitive Science, Psychology Press, New York, 1984, pp. 21–49. Ch. 2.
- [13] M.T. Maybury (Ed.), New Directions in Question Answering, AAAI Ser., MIT Press, Cambridge MA, 2004.
- [14] A. Mishra, S.K. Jain, A survey on question answering systems with classification, J. King Saud Univ. - Comput. Inf. Sci. 28 (2016) 345–361.
- [15] R.J. Nelson, Machine models for cognitive science, Philos. Sci. 54 (1987) 391–408.
- [16] M. Nivat, 40th anniversary of EATCS, Bull. Eur. Assoc. Theor. Comput. Sci. EATCS 108 (October 2012) 145–147.
- [17] N. Pippenger, On simultaneous resource bounds, in: 20th Symp. Foundations of Computer Science, Proceedings, IEEE Press, 1979, pp. 307–311.
- [18] Z.W. Pylyshyn, Computing in cognitive science, in: M. I. Posner (Ed.), Foundations of Cognitive Science, MIT Press, Cambridge, MA, 1989, pp. 49–91.
- [19] J. van Leeuwen, J. Wiedermann, Beyond the Turing limit: evolving interactive systems, in: L. Pacholski, P. Ružička (Eds.), Theory and Practice of Informatics, Proc. 28th Conference, SOFSEM 2001, in: Lecture Notes in Comput. Sci., vol. 2234, Springer-Verlag, Berlin, 2001, pp. 90–109.
- [20] H. Vollmer, Introduction to Circuit Complexity – A Uniform Approach, Springer-Verlag, Berlin, 1999.
- [21] J. Wiedermann, The computational limits to the cognitive power of the neoroidal tabula rasa, J. Exp. Theor. Artif. Intell. 15 (3) (2003) 267–279.