

facebook

Artificial Tasks for Artificial Intelligence

Antoine Bordes + Jason Weston, Sumit Chopra, Tomas Mikolov, Armand Joulin, Sasha Rush & Léon Bottou
Facebook AI Research
ICLR – San Diego – May 7, 2015

Menu

1. How building intelligent machines *“are we there yet?”*
2. A look back *“a blast from the past”*
3. Artificial tasks for learning for AI *“a virtuous circle”*
4. A look forward *“is virtuous also incestuous?”*
5. Wrap up *“a leap of faith?”*

Building intelligent machines

An era of success

- Many breakthroughs recently

- Object detection (Krizhevsky et al. 13)
- Speech recognition (Hinton et al. 12)
- Word embeddings (Collobert et al. 11) (Mikolov et al. 13)
- Machine translation (Sutskever et al. 14)



- Ingredients

1. Models with high capacity and representation power (CNNs, RNNs, LSTMs, etc.)
2. Lots of (supervised) data

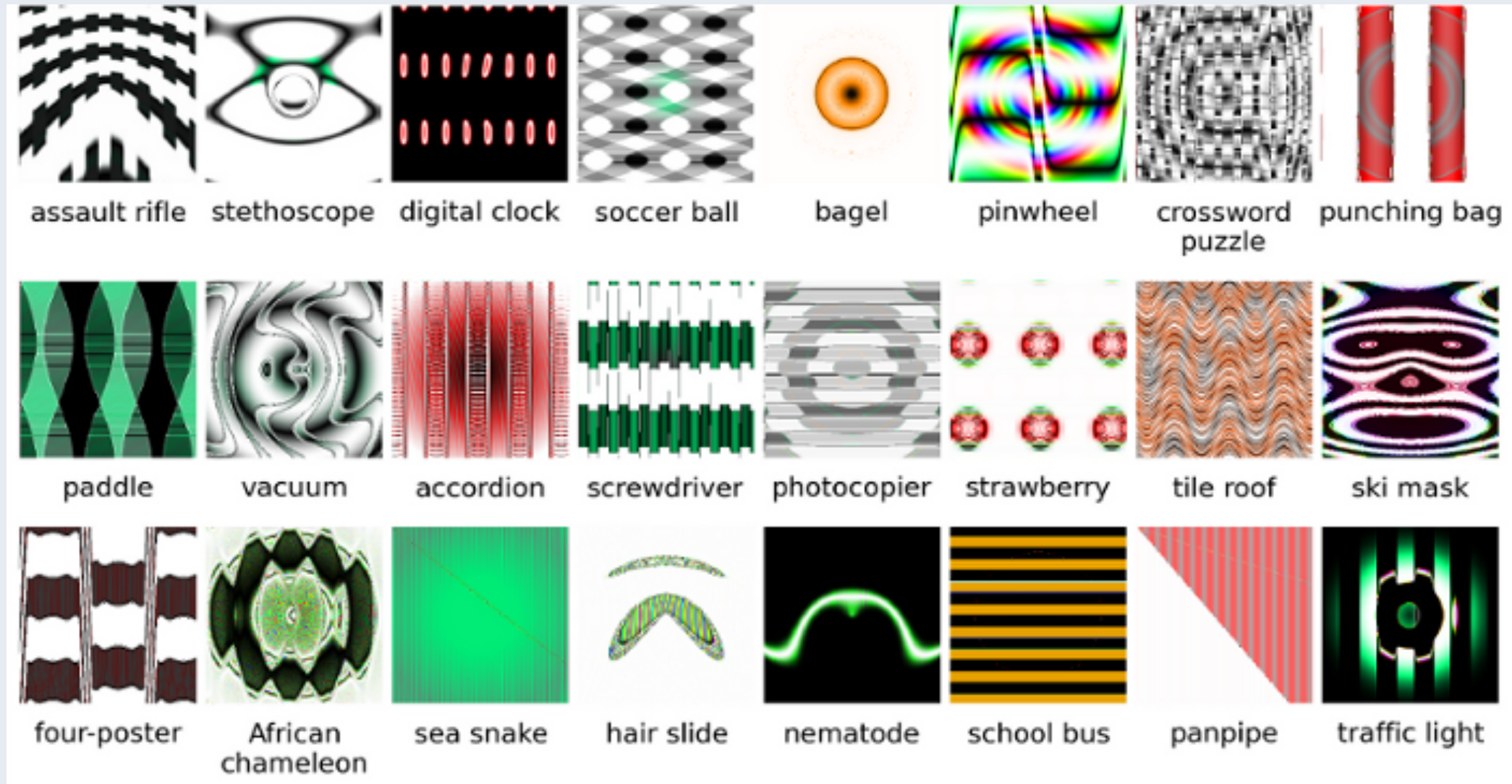
Except for word embeddings..

So AI is solved?

- Not quite.
- Success of Big data statistics cleverly used.
 - Labeled data is crucial.
 - So are fast and powerful computers (m/billions of examples/parameters).
 - Models and learning algorithms are not specifically new.

- Reasoning is still limited.
- “Concept \neq statistics” (Bottou 15)

Convnets can be fooled



(Nguyen et al. CVPR15)

Caption generation

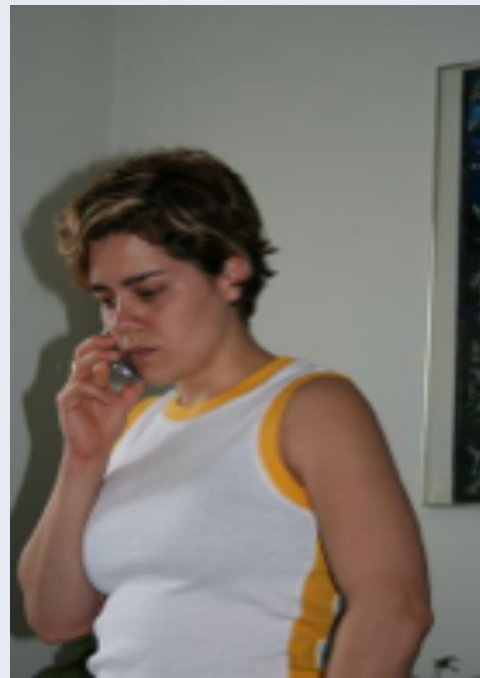
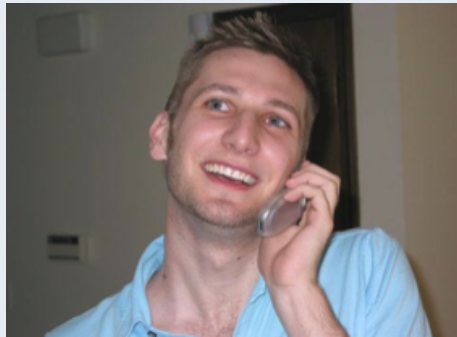
A herd of giraffes walk down the street in the middle of some trees.



From a talk of
(Zitnick 14)

Action (or multiple objects ?) detection

Detection the **action** “phoning”



(Oquab et al. CVPR 2014)

MT still does not fully understand ...

Ensemble of 8 LSTMs + unknown-words (Luong et al. ACL15):

[eng] But concerns have grown after Mr Mazanga was quoted as saying Renamo was abandoning the 1992 peace accord.

[*eng*] But concerns have grown after Mr Mazanga declared that Renamo was the 1992 peace accord.

... nor do Q&A systems.

Embedding-based model of (Bordes et al. EMNLP14):

- What is Jimi Hendrix Purple Haze about?

A:guitar

A:drug_overdose

- What country was Slovakia?

A:austria, A:czech_republic

A:poland, A:czechoslovakia, A:france, A:hungary,

Big Data => Big AI ?

- Can we solve AI with current models with more data and more computing power ?
 - We can certainly improve (a lot) on many (well defined) tasks.
 - Training data will never cover the whole range of possible situations
 - It is always a proxy (Imagenet is a proxy for vision)
 - What happens when train/test distributions drift?

Example

How can we learn a conversational agent? The system must:

- Adapt quickly to the context: target distribution evolves
- Learn incrementally and build new knowledge from heterogeneous sources
- Have a long-term structured memory

Limits of Big Model + Big Data

- Training and evaluation on real large-scale data is difficult:
 - Real large-scale data is complex, noisy, unlabeled... → big infrastructure
 - Interpretation of success or failure is complex
- Complicates the design of innovative learning systems

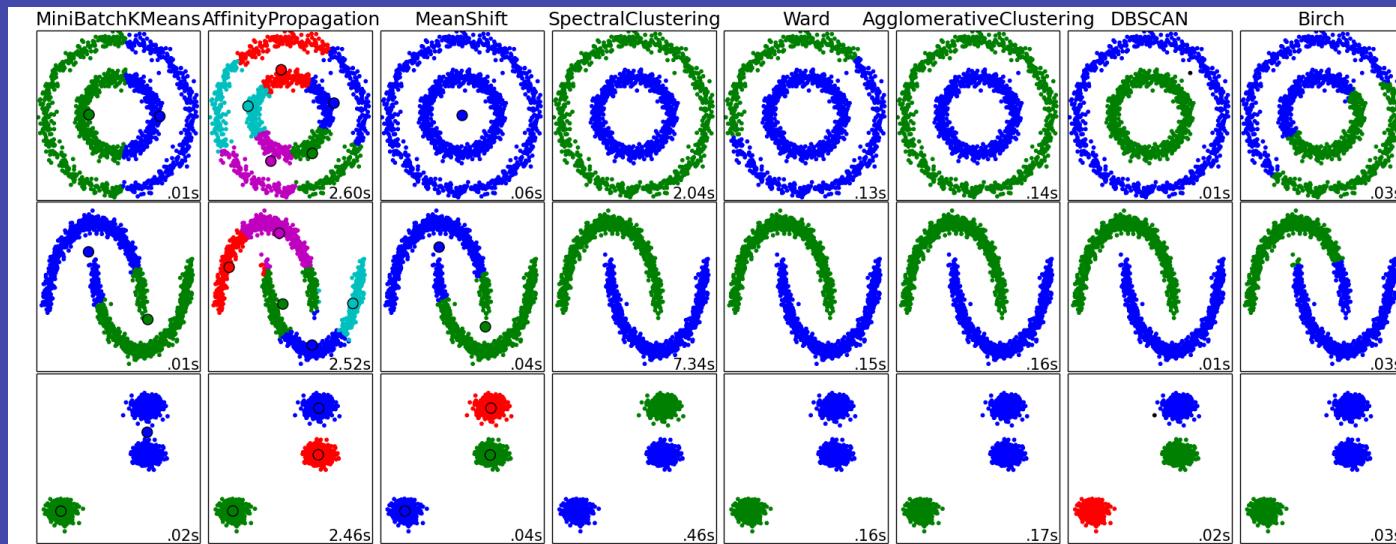
This talk: Artificial environments for training algorithms.

- Total control on the complexity of the tasks/reasoning
- Clear interpretation of results
- Assessment that the system behaves the way we want
- Challenge: how transfer from artificial to real conditions?

Looking backwards

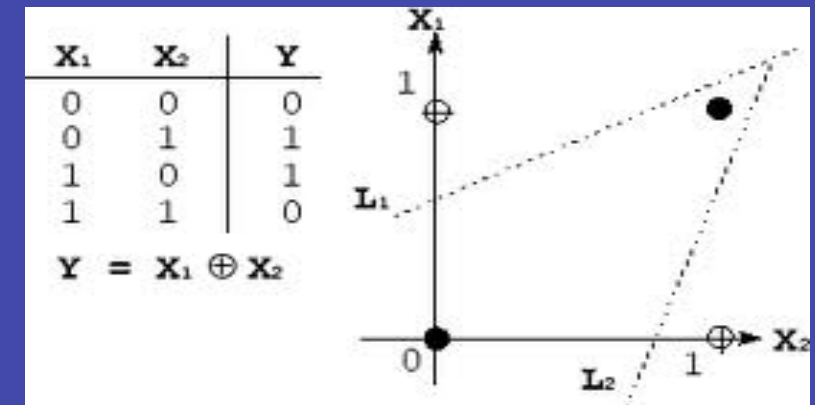
Artificial problems in ML

Two moons and friends (clustering)



Toy/artificial problems in ML: crucial for demonstrating and assessing the usefulness/efficiency of new algorithms

XOR (neural networks)

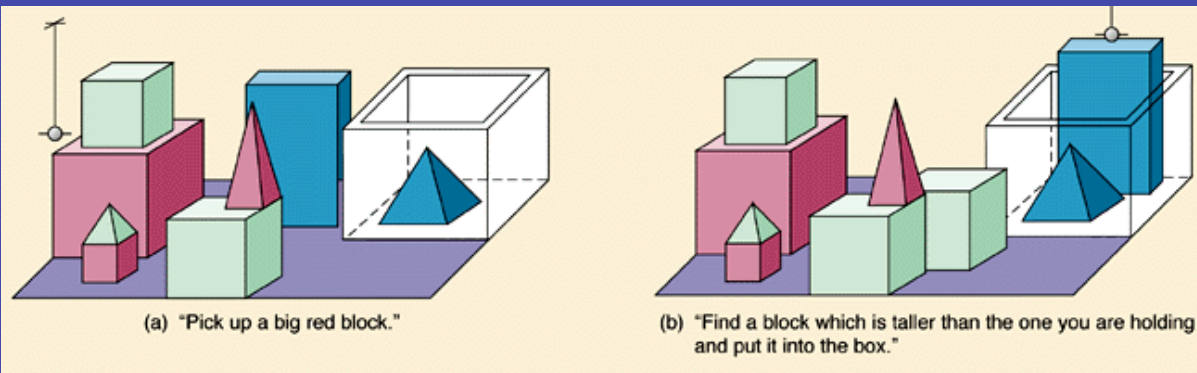


Many in the UCI repository
(regression, classification)



And in AI?

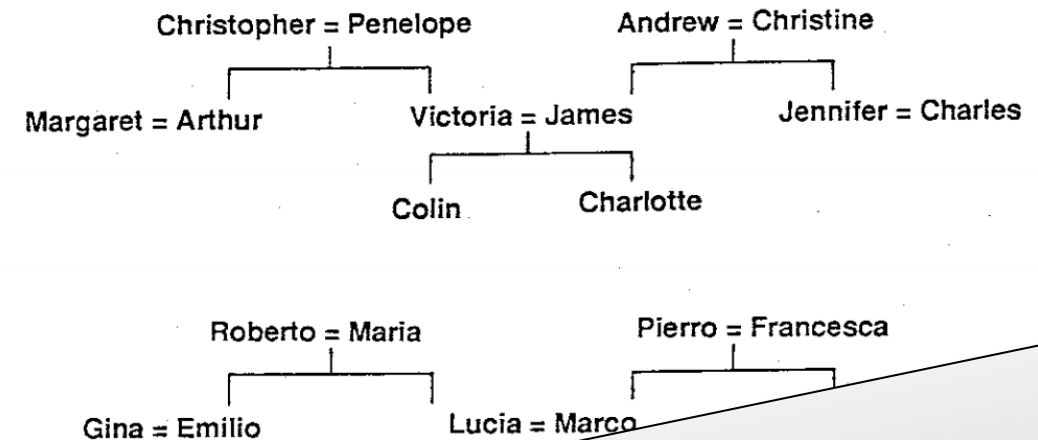
Block Worlds (Winograd 71)



Person: PICK UP A BIG RED BLOCK.
Computer: OK. (does it)
Person: GRASP THE PYRAMID.
Computer: I DON'T UNDERSTAND WHICH PYRAMID YOU MEAN.
Person: FIND A BLOCK WHICH IS TALLER THAN THE ONE YOU
Computer: BY "IT", I ASSUME YOU MEAN THE BLUE BLOCK.
Computer: OK. (does it)
Person: WHAT DOES THE BOX
Computer: THE BOX IS A BOX.
Person: STACK UP TWO PYRAMIDS.

- SHRLDU completely scripted
- Impossible to scale up

Family trees (Hinton CogSci86)



- Emphasis on learning
- Still used as benchmark
- Recent (large-scale) successors
(Bordes et al. AAAI11) (Dong et al. KDD14)

Artificial tasks for learning for AI

Evaluating and learning systems for AI

- Recent effort to build controlled evaluation environments:
 - Project ARISTO (Allen Institute for AI): pass 5th grader science exams
 - Winograd Schema Challenge (Levesque AAAI11): ~150 schemas

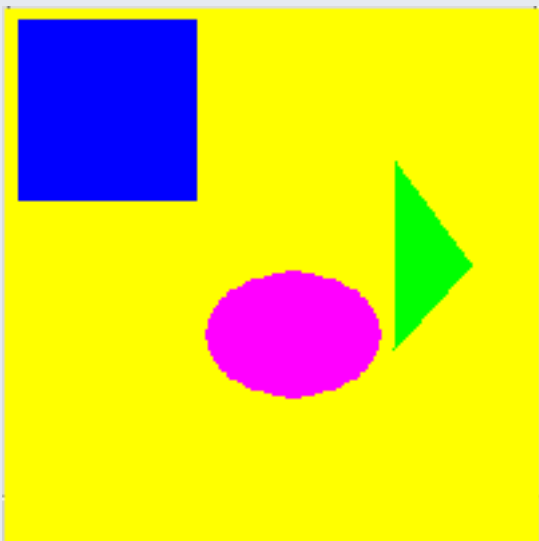
The city councilmen refused the demonstrators a permit because they [feared/advocated] violence. Who [feared/advocated] violence? **A: The city councilmen/the demonstrators**

The trophy doesn't fit into the brown suitcase because it's too [small/large]. What is too [small/large]? **A: The suitcase/the trophy**

If we want to motivate the creation of new learning algorithms, training conditions are crucial and should be controlled too.

ShapeseT

- A rebirth of **block worlds**: less ambitious but **emphasis on learning**.



TASK	Example of question	Answer
Color	There is a small triangle. What color is it?	Green
Shape	What is the shape of the green object?	Triangle
Size	There is a triangle on the right. Is it rather small or bigger?	Small
Size (relative)	There is a square on the top left. Is it smaller or bigger than the triangle?	Bigger
Location	Is the blue square at the top or at the bottom?	At the top
Location (relative)	There is a purple ellipse. Is it on the right or on the left of the triangle?	On the left

- Learn to answer questions given images generated from a simulation.

(Breuleux et al. 08)

See also the Synthetic Visual Reasoning Test (Fleuret & Geman 10)

Sequences

- A range of various basic tasks (usually) at the character level:

- Copying.

- **Sorting.**

- Associative recall.

- Dynamic n-grams.

- Counting:

Sequence generator	Example
$\{a^n b^n \mid n > 0\}$	aab ba aab bb ab aaaa ab bbbb
$\{a^n b^n c^n \mid n > 0\}$	aab cc ab ca aaaab bbbb cccc
$\{a^n b^n c^n d^n \mid n > 0\}$	aab dd aaab bb cc ddd abcd
$\{a^n b^{3n} \mid n > 0\}$	aab bb baaab bbbb bbbab bb
$\{a^n b^n \mid n > 0\}$	aab ca aabb cccc abcc
$\{a^n b^n \mid n > 0\}$	aab ca cc bcbbbbcba bc bbb
$\{a^n b^n \mid n > 0\}$	aab ca aabb cccc abc

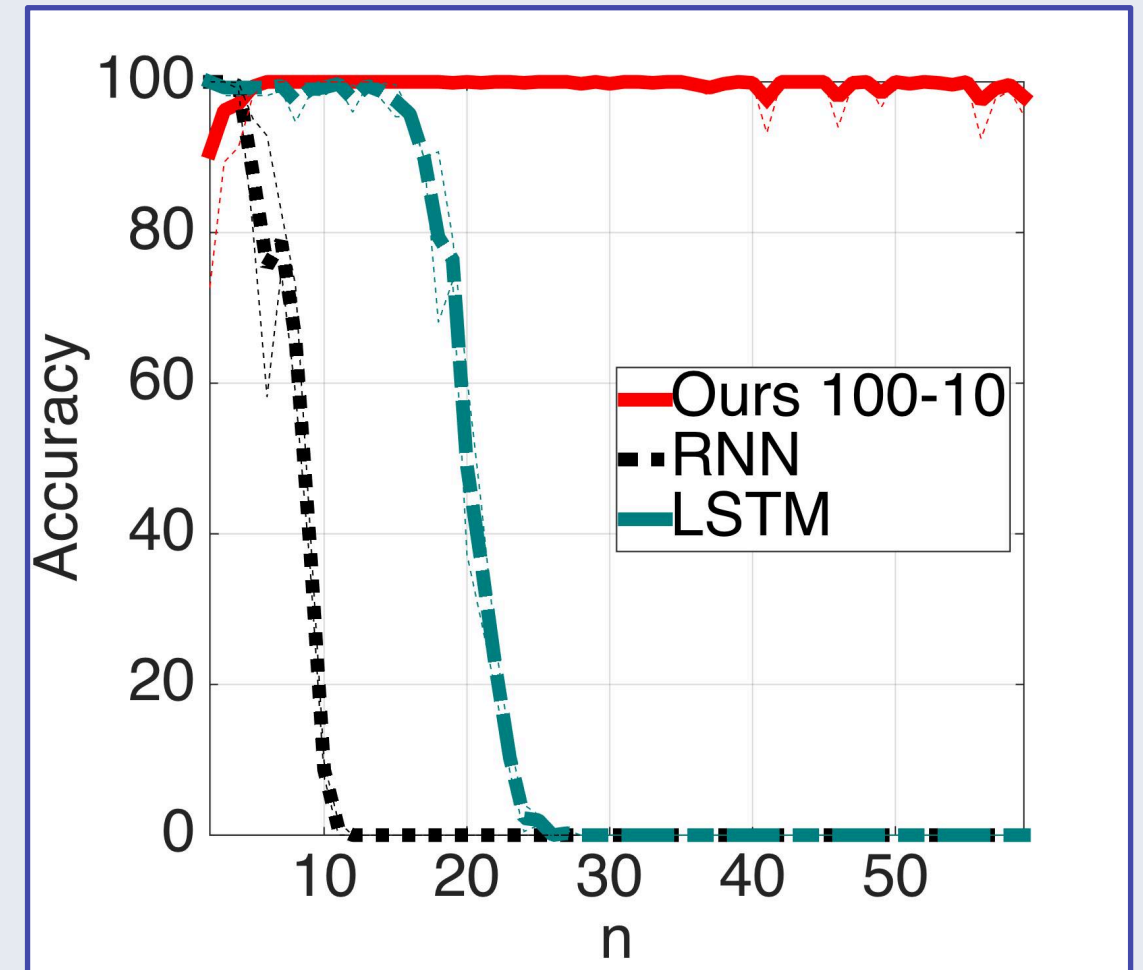
- Test the ability to generalize to long-term dependencies.
- Goal: generate sequences never seen in training.

- Recently already proven useful:
- Neural Turing Machines (Graves et al. 14)
- Stack-augmented RNNs (Joulin & Mikolov 15)

Stack-augmented RNNs

Task $a^n b^{2n}$

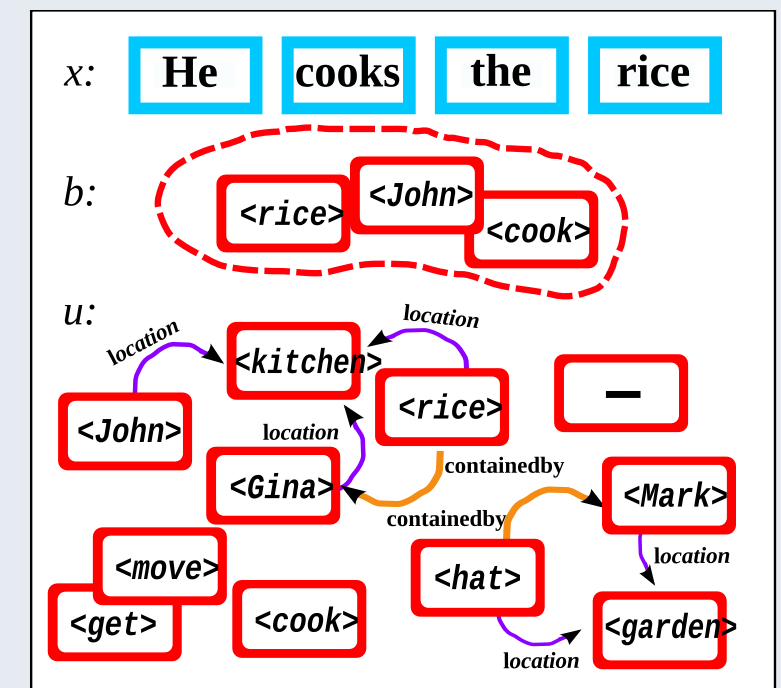
current	next	prediction	proba(next)	action		stack1[top]	stack2[top]
b	a	a	0.99	POP	POP	-1	0.53
a	a	a	0.99	PUSH	POP	0.01	0.97
a	a	a	0.95	PUSH	PUSH	0.18	0.99
a	a	a	0.93	PUSH	PUSH	0.32	0.98
a	a	a	0.91	PUSH	PUSH	0.40	0.97
a	a	a	0.90	PUSH	PUSH	0.46	0.97
a	b	a	0.10	PUSH	PUSH	0.52	0.97
b	b	b	0.99	PUSH	PUSH	0.57	0.97
b	b	b	1.00	POP	PUSH	0.52	0.56
b	b	b	1.00	POP	PUSH	0.46	0.01
b	b	b	1.00	POP	PUSH	0.40	0.00
b	b	b	1.00	POP	PUSH	0.32	0.00
b	b	b	1.00	POP	PUSH	0.18	0.00
b	b	b	0.99	POP	PUSH	0.01	0.00
b	b	b	0.99	POP	POP	-1	0.00
b	b	b	0.99	POP	POP	-1	0.00
b	b	b	0.99	POP	POP	-1	0.01
b	a	a	0.99	POP	POP	-1	0.56



(Das et al. CogSci92) (Joulin & Mikolov 15)

Text adventure games

- We explore such games to create learning environments
- A simulated world, like a text adventure game, can generate stories
 - From actions, sentences are produced using a simple grammar
 - This allows to ground language into actions
 - Difficulty/complexity is controlled
 - Training and evaluation data are provided
 - Evaluation through question answering is easy



(Bordes et al. AISTATS10)

Simulation commands

- go <place>
- get <object>
- get <object1> from <object2>
- put <object1> in/on <object2>
- give <object> to <person>
- drop <object>
- look
- inventory
- examine <object>

+ 2 commands for "gods" (superusers):

- create <object>
- set <obj1> <relation> <obj2>

Example

Simple grammar

Command format

```
jason go kitchen  
jason get milk  
jason go office  
jason drop milk  
jason go bathroom  
where is milk ?    A: office  
where is jason? A: bathroom
```

Story

Jason went to the kitchen.
Jason picked up the milk.
Jason travelled to the office.
Jason left the milk there.
Jason went to the bathroom.
Where is the milk now? **A: office**
Where is Jason? **A: bathroom**

A collection of tasks

- We created 20 tasks:
 - Paper: (Weston et al. 15) arxiv.org/abs/1502.05698
 - Data: facebook.ai/babi
- Each task checks one skill that a reasoning system should have.
- We look for systems able to solve all tasks: *no task specific engineering*.

We postulate that *performing well on all of them is a pre-requisite for any system aiming at understanding language and able to reason.*

(T1) Single supporting fact “where is actor”

- Questions where a single supporting fact, previously given, provides the answer.
- Simplest case of this: asking for the location of a person.

John is in the playground.
Bob is in the office.
Where is John? A:playground

SUPPORTING FACT



(T2) Two supporting facts “where is actor+object”

- Harder task: questions where **two supporting statements** have to be **chained** to answer the question.



John is in the playground.

Bob is in the office.

John picked up the football.

Bob went to the kitchen.

Where is the football? **A:playground**

Where was Bob before the kitchen? **A:office**

SUPPORTING FACT

SUPPORTING FACT

- To answer the first question *Where is the football?* both John picked up the football and John is in the playground are supporting facts

(T3) Three supporting facts

- Similarly, one can make a task with **three supporting facts**:

John picked up the apple.
John went to the office.
John went to the kitchen.
John dropped the apple.
Where was the apple before the kitchen? **A:office**



- The first three statements are all required to answer this.

(T4) Two argument relations: subj vs. obj.

- To answer questions the **ability to differentiate and recognize subjects and objects is crucial.**
- We consider the extreme case: **sentences feature re-ordered words:**

The office is north of the bedroom.
The bedroom is north of the bathroom.
What is north of the bedroom? **A:office**
What is the bedroom north of? **A:bathroom**

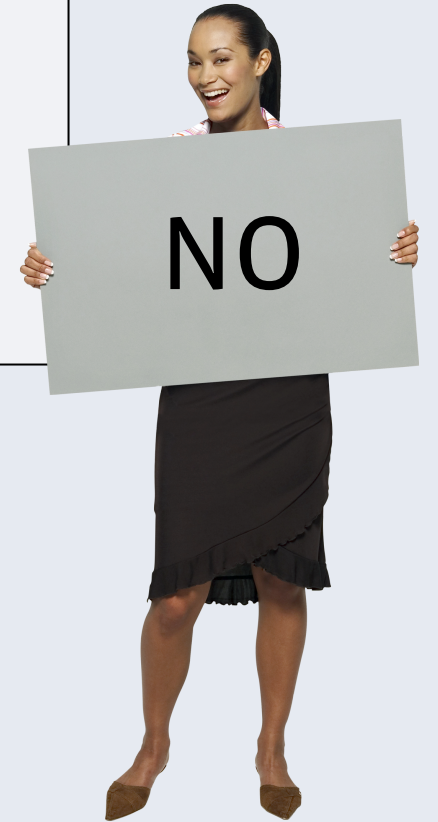


- The two questions above have exactly the same words, but in a different order, and different answers.
- So **a bag-of-words will not work.**

(T6) Yes/No questions

- This task tests, in the simplest case possible (with a single supporting fact) the ability of a model to answer true/false type questions:

John is in the playground.
Daniel picks up the milk.
Is John in the classroom? A:no
Does Daniel have the milk? A:yes



(T7) Counting

- This task tests the ability of the QA system to perform **simple counting operations**, by asking about the number of objects with a certain property:

Daniel picked up the football.

Daniel dropped the football.

Daniel got the milk.

Daniel took the apple.

How many objects is Daniel holding? **A:two**



(T17) Positional reasoning

- This task tests **spatial reasoning**, one of many components of the classical block world :

The triangle is to the right of the blue square.
The red square is on top of the blue square.
The red sphere is to the right of the blue square.
Is the red sphere to the right of the blue square? **A:yes**
Is the red square to the left of the triangle? **A:yes**



- Close from Shapeset or block worlds, **with no vision input**.
- The Yes/No task (6) is a prerequisite.

(T18) Reasoning about size

- This task requires reasoning about relative size of objects :

The football fits in the suitcase.
The suitcase fits in the cupboard.
The box of chocolates is smaller than the football.
Will the box of chocolates fit in the suitcase? A:yes



- Inspired by the commonsense reasoning examples of the Winograd schema challenge (Levesque AAAI11)
- Tasks 3 (three supporting facts) and 6 (Yes/No) are prerequisites.

(T19) Path finding

- In this task the goal is to **find the path** between locations:

The kitchen is north of the hallway.

The den is east of the hallway.

How do you go from den to kitchen? **A:west,north**



- This task is difficult because it effectively **involves search**.

Training on 1k stories

Dashboard

Weak supervised

Fully supervised

TASK	N-grams	LSTMs	StructSVM + COREF + SRL	Memory Networks
T1. Single supporting fact	36	50	PASS	PASS
T2. Two supporting facts	2	20	74	PASS
T3. Three supporting facts	7	20	17	PASS
T4. Two arguments relations	50	61	PASS	PASS
T5. Three arguments relations	20	70	83	84
T6. Yes/no questions	49	48	PASS	49
T7. Counting	52	49	69	73
T8. Sets	40	45	70	87
T9. Simple negation	62	64	PASS	62
T10. Indefinite knowledge	45	44	PASS	50
T11. Basic coreference	29	72	PASS	PASS
T12. Conjunction	9	74	PASS	PASS
T13. Compound coreference	26	PASS	PASS	PASS
T14. Time reasoning	19	27	PASS	PASS
T15. Basic deduction	20	21	PASS	PASS
T16. Basic induction	43	23	24	PASS
T17. Positional reasoning	46	51	61	48
T18. Size reasoning	52	52	62	68
T19. Path finding	0	8	49	4
T20. Agent's motivation	76	91	PASS	PASS

Looking forward

How transfer from artificial to real data?

Tasks are useful on their own as pre-requisite tests for reasoning.

To (eventually) scale up to real language:

1. No model should be tailored for a task alone, nor for the tasks only.
2. We should look for models able to learn incrementally faster new tasks.
3. We have 20 AI tasks. We will create others : not a definitive set!
 - The simulation is parameterized to ramp up the complexity
 - Annotators could be used to generate real language from it

Simulation control panel



Symbols

- Can the system switch to other languages?
- And other (simpler) symbolic systems?



Memory

- How far should one remember?
- Is an external source of knowledge necessary?



Linguistics

- How is reasoning altered by ambiguities?
- And by embedded clauses?



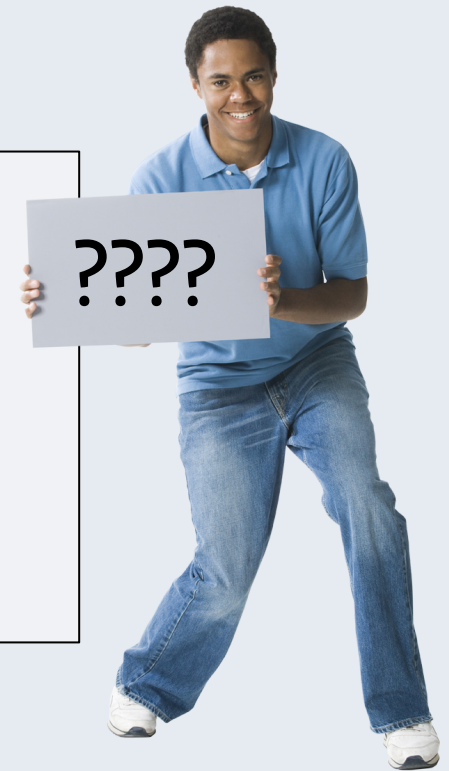
Reasoning

- How many facts should be chained together?
- How many examples does the system need?

Playing with the **symbols** knob

- We can use other languages and produce equivalent datasets. This is what Task (T13) looks like in Hindi:

sita aur badri galiyarey mein chale gaye
uske upraant wo daftar mein chale gaye
priya aur mohit daftar mein chale gaye
uske baad wo galiyarey mein chale gaye
badri is samay kahan hai ? **A:daftar**



- How would perform the SVM in this setting where coref/SRL might be worse?

Playing with the **symbols** knob

- We can also shuffle the letters and produce other equivalent datasets:

Sbdm ip im vdu yonrckblms.

Abf ip im vdu bhhigu.

Sbdm yigaus ly vdu hbbvfnoo.

Abf zumv vb vdu aivgdum.

Mduku ip vdu hbbvfnoo? **A:yonrckblms**

Mduku znp Abf fuhbku vdu aivgdum? **A:bhhigu**



- The reasoning is still learnable but usual NLP systems can not access it.

Playing with the **memory** knob

- We can tune the **distance between supporting facts and the question, with irrelevant facts**. Hence, Task (T1) can become:

John is in the playground.
Bob is in the office.
Ringo went to San Diego.
Paul attended ICLR.
George played the guitar.
Ringo bought drums.
They jumped in a yellow submarine.
And they flew in the sky with Diamond.
....
....
Where is John? **A:playground**



Playing with the **memory** knob

- We can also require the system to learn to use external resources to be able to solve the task (common-sense):

John went to the restaurant.
John ordered a burger.
John left a big tip.
Did John like the restaurant? **A:yes**
Is John a vegan? **B:no**



- Here: all information needed to answer is not only in the training stories.

Playing with the **linguistics** knob

- Task (T20) tests **the simplest type of coreference**, that of detecting the nearest referent, for example:

Daniel was in the kitchen.
Then he went to the studio.
Sandra was in the office.
Where is Daniel? **A:studio**

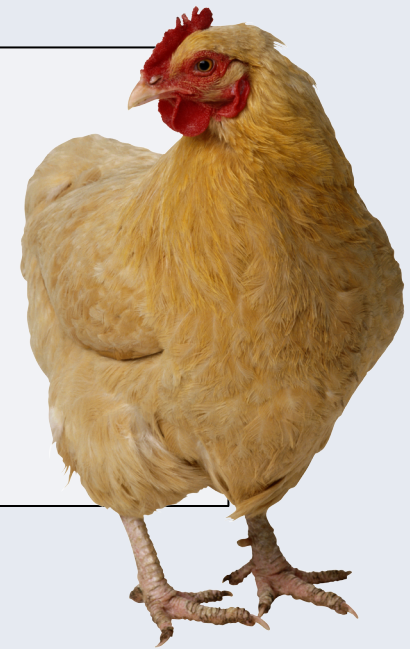


- Increasing difficulty:
 - + **flip order of last two statements.**
 - +++ adapt a real coreference dataset into a Q&A format.

Playing with the **linguistics** knob

- Task (T13) tests coreference when **the pronoun can refer to multiple actors**:

Daniel and Sandra journeyed to the office.
Then they went to the garden.
Sandra and John travelled to the kitchen.
After that they moved to the hallway.
Where is Daniel? **A:garden**



Playing with the **linguistics** knob

- Task (T14) tests **understanding the use of time expressions**:

In the afternoon Julie went to the park.
Yesterday Julie was at school.
Julie went to the cinema this evening.
Where did Julie go after the park? **A:cinema**



- Much harder difficulty**: adapt a real time expression labeling dataset into a question answer format, e.g. (Uzzaman et al. 12).

Playing with the reasoning knob

- Task (T8) tests the ability to produce a set of single word:

Daniel picks up the football.
Daniel drops the newspaper.
Daniel picks up the milk.
What is Daniel holding? **A:milk,football**



- The task above can be seen as a QA task related to database search. We could also consider the following question types:
 - Intersection: *Who is in the park carrying food?*
 - Union: *Who has milk or cookies?*
 - Set difference: *Who is in the park apart from Bill?*

Playing with the reasoning knob

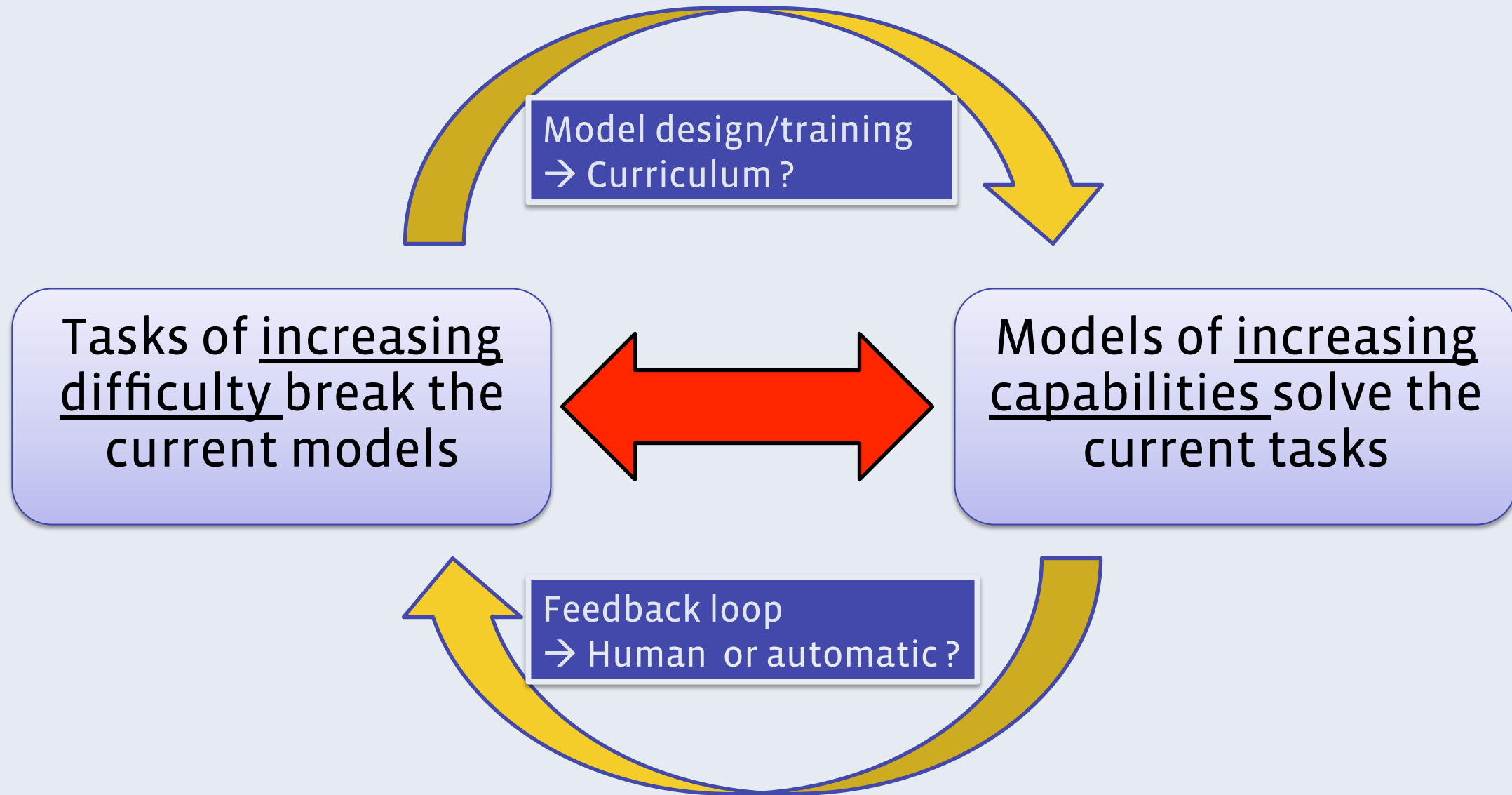
- We can also generate stories with more complex underlying rules:

Daniel picked eight bananas.
Daniel gave a quarter of these bananas to Paul.
Paul ate half of his bananas.
Paul bought more bananas to triple his stock.
How many bananas does Paul have? A: three



- Might be careful to decorrelate this from linguistics difficulties.

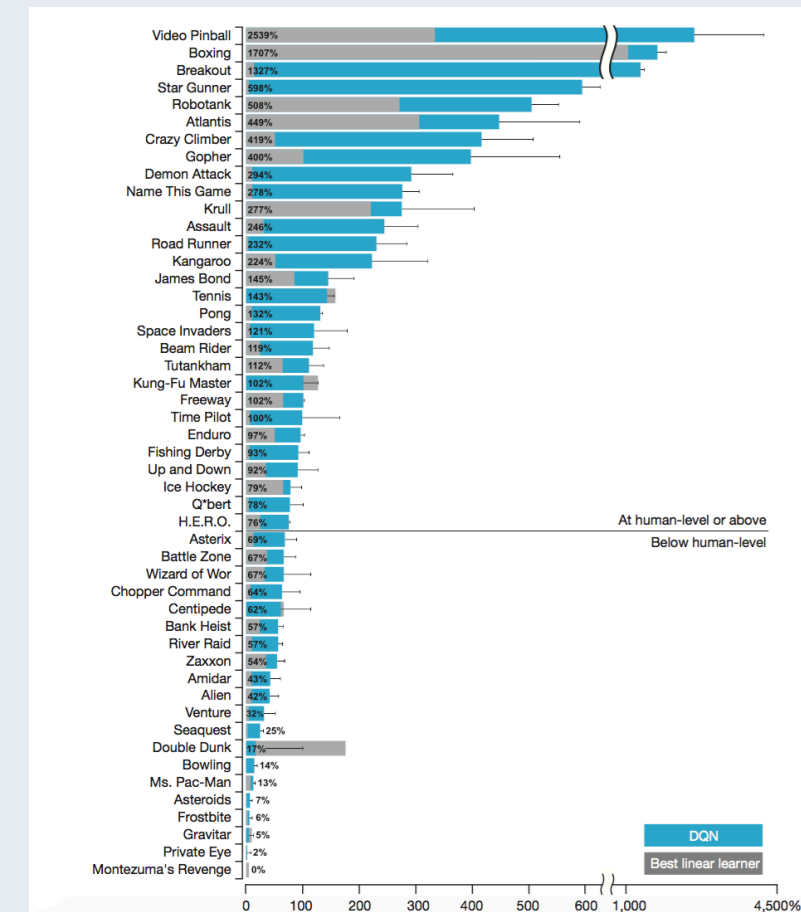
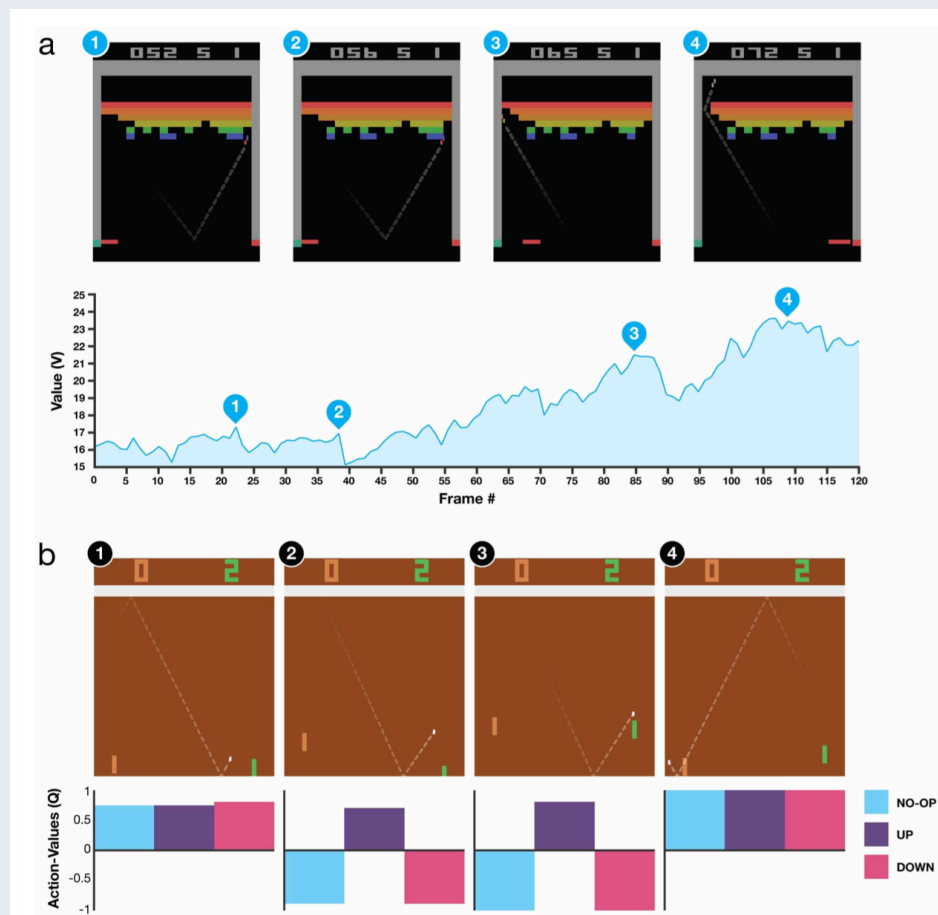
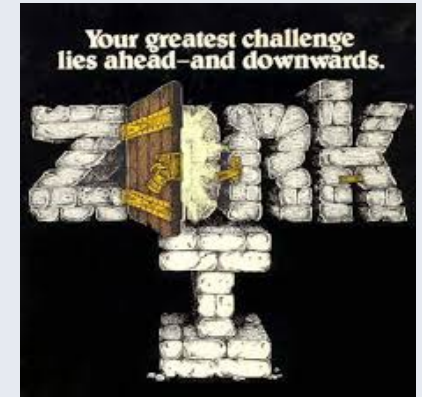
A virtuous circle



Beware! The circle should be virtuous but not incestuous

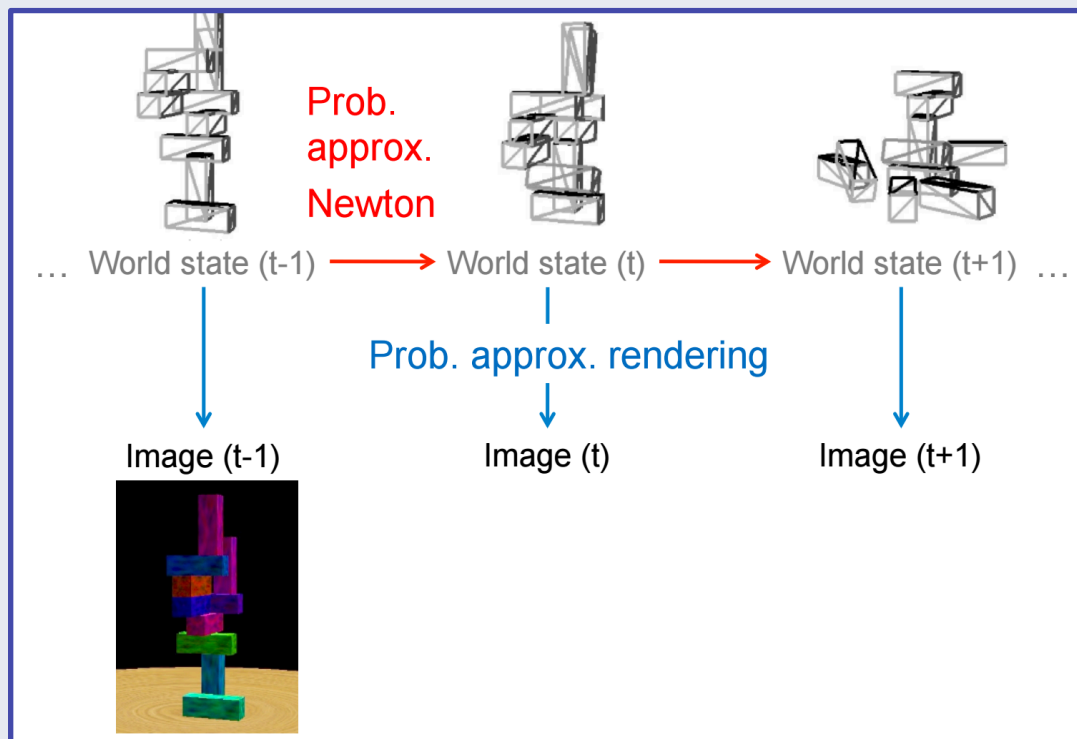
Computer games

- Our simulation is based on old-style text adventure games
 - Old games offer a great variety of controlled environments
- Atari games for reinforcement learning (Mnih et al. Nature15)

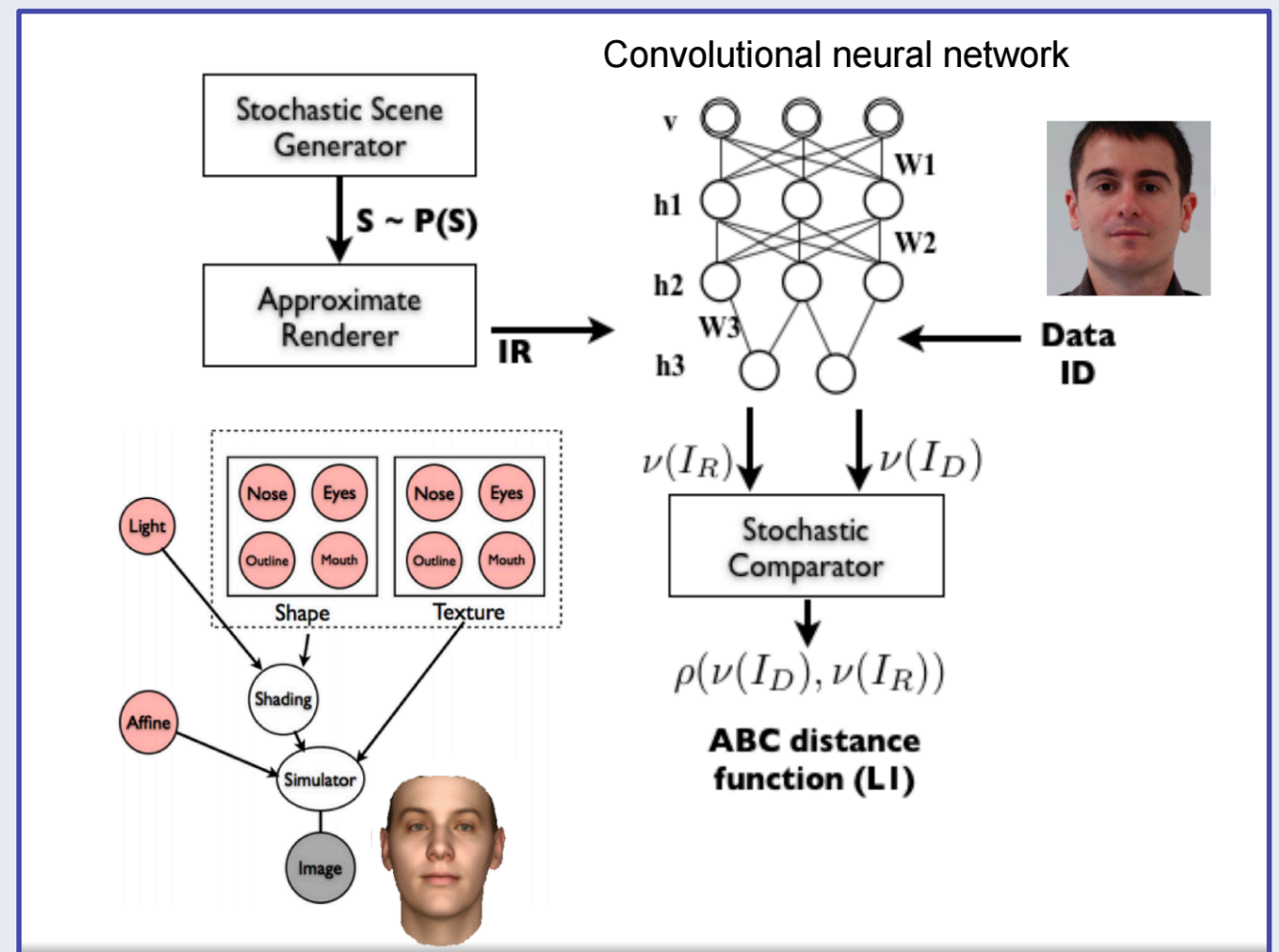


A game engine within the model?

- Prior knowledge about the structure of the world is pre-wired in the model using a game engine



(Taglia et al. PNAS13)

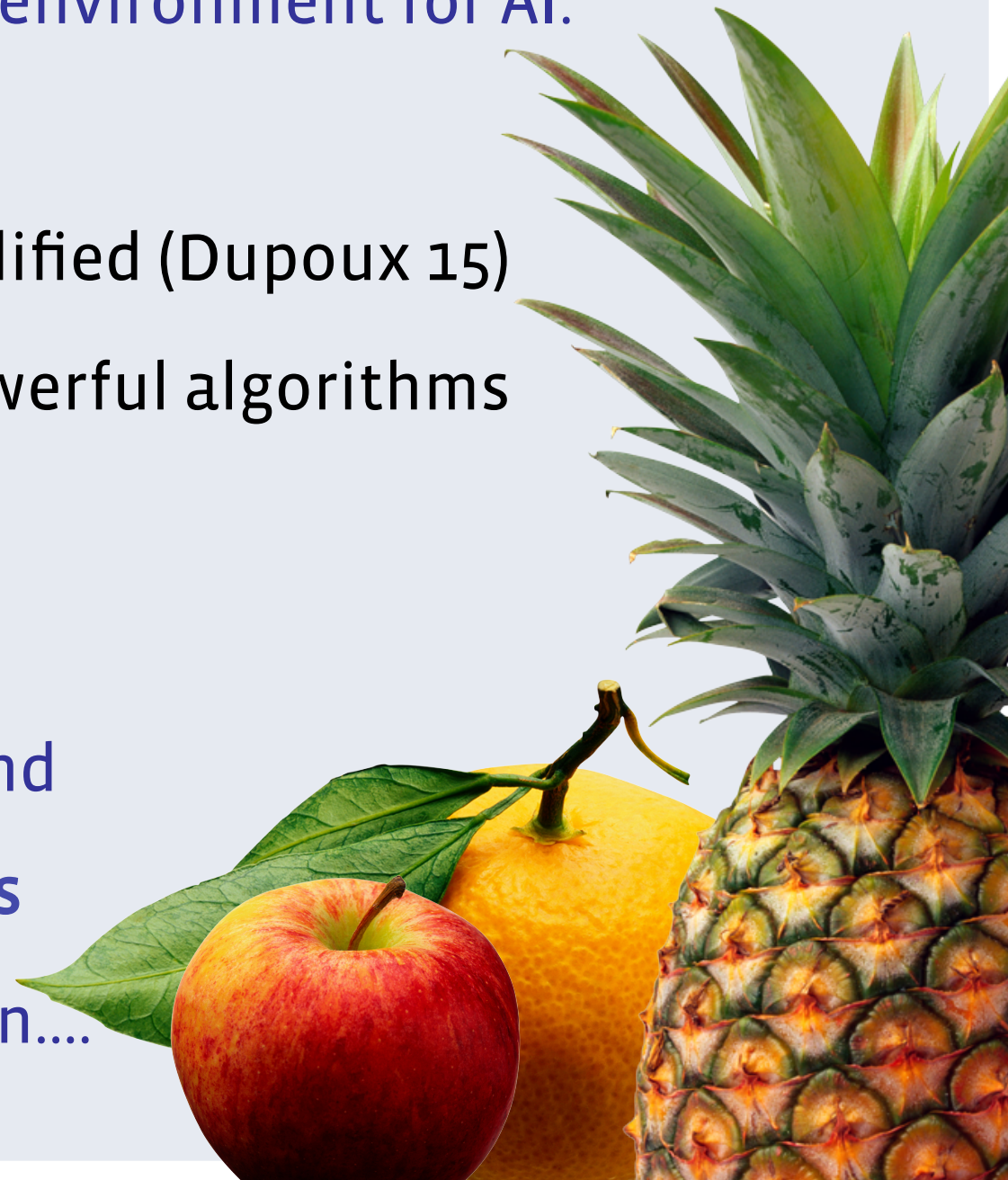


(Kulkarni et al. CVPR15)

Wrap up

On the need of artificial tasks

- We need a controlled training and testing environment for AI.
- Artificial tasks of increasing difficulty.
- Their design can be modified or even simplified (Dupoux 15)
- Drives the design of increasingly more powerful algorithms
- Our hope is that a feedback loop of:
 1. Developing tasks that break models, and
 2. Developing models that can solve tasks... leads in a fruitful research direction....



facebook

(c) 2009 Facebook, Inc. or its licensors. "Facebook" is a registered trademark of Facebook, Inc.. All rights reserved. 1.0