

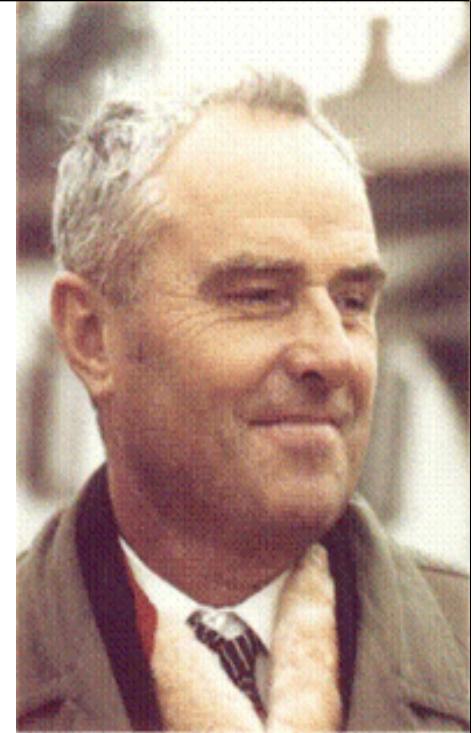
Unsupervised Minimax: nets that fight each other

Jürgen Schmidhuber
The Swiss AI Lab IDSIA
Univ. Lugano & SUPSI
<http://www.idsia.ch/~juergen>

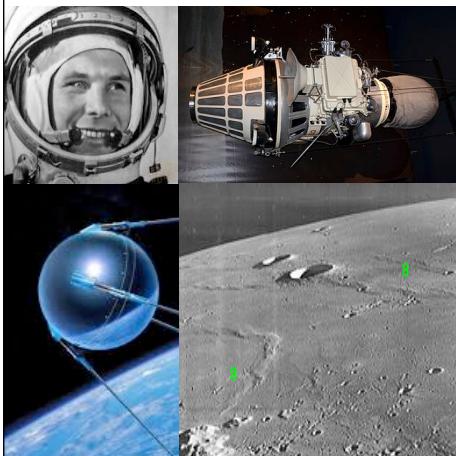
NNAISENSE

Jürgen Schmidhuber
You_again Shmidhoobuh

Supervised Deep Learning in feedforward networks:
A. G. Ivakhnenko & Lapa (Ukraine, since 1965). Deep nets with arbitrary number of layers & polynomial activation functions learn internal representations. Layer-wise training by regression analysis: learn numbers of layers and units per layer. Prune superfluous units



O.Г. Івахненко (1967 р.)



USSR around this time: start of space age (Sputnik 1957), first man in space (Gagarin, 1961), biggest bomb ever (Tsar Bomba 1961), first robots on the moon (Luna 9, 1966) and another planet (Venera 7, 1970)



USSR also was home to many of the greatest mathematicians

Deep nets with 8 layers already back in 1970. Still used in the 2000s

“Modern” Backpropagation (BP, 1970)

<http://people.idsia.ch/~juergen/who-invented-backpropagation.html>

Continuous BP in Euler-LaGrange Calculus + Dynamic Programming: Kelley 1960, Bryson 1961. BP through chain rule only: Dreyfus 1962. ‘**Modern BP**’ or **automatic differentiation (AD) in sparse, discrete, NN-like nets**: **Linnainmaa 1970**. Weight changes: Dreyfus 1973. BP applied to NNs: Werbos 1982 (first thoughts: 1974). Experiments with 1000 times faster computers yield useful internal representations: Rumelhart et al 86.

Recurrent NNs: e.g., Williams, Werbos, Robinson, 1980s...



General Purpose Deep Learning with RNNs since 1991

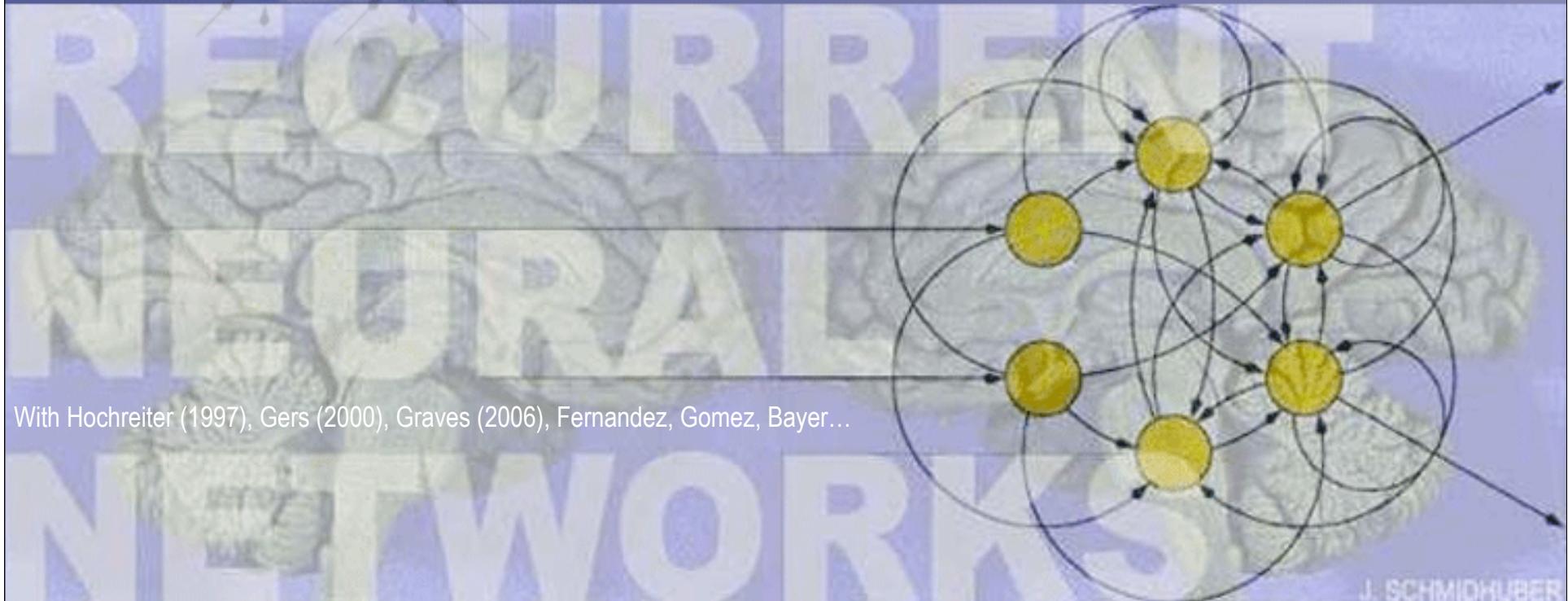
<http://www.idsia.ch/~juergen/firstdeeplearn.html>

Neural history compressor:
unsupervised pre-training of RNN
stack or hierarchy through predictive
coding. Compress chunker RNN
(teacher) into automatizer RNN
(student) also re-trained on previous
skills. Experiments: depth >1000

LONG SHORT-TERM MEMORY

But then supervised RNNs took over, through LSTM.

Now in the AI on your phone

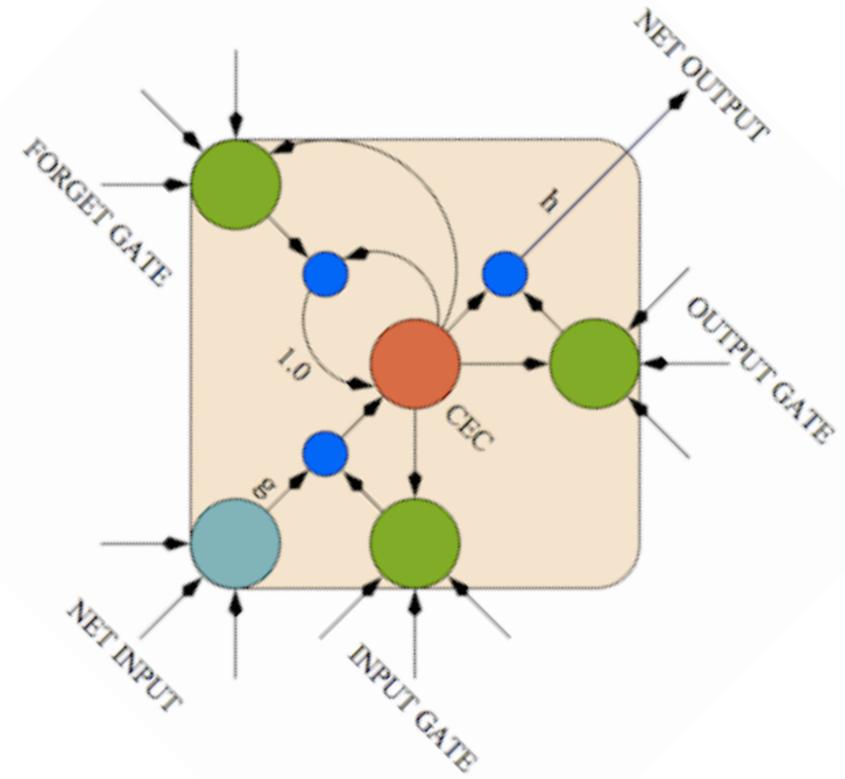


Today's LSTM shaped by my:

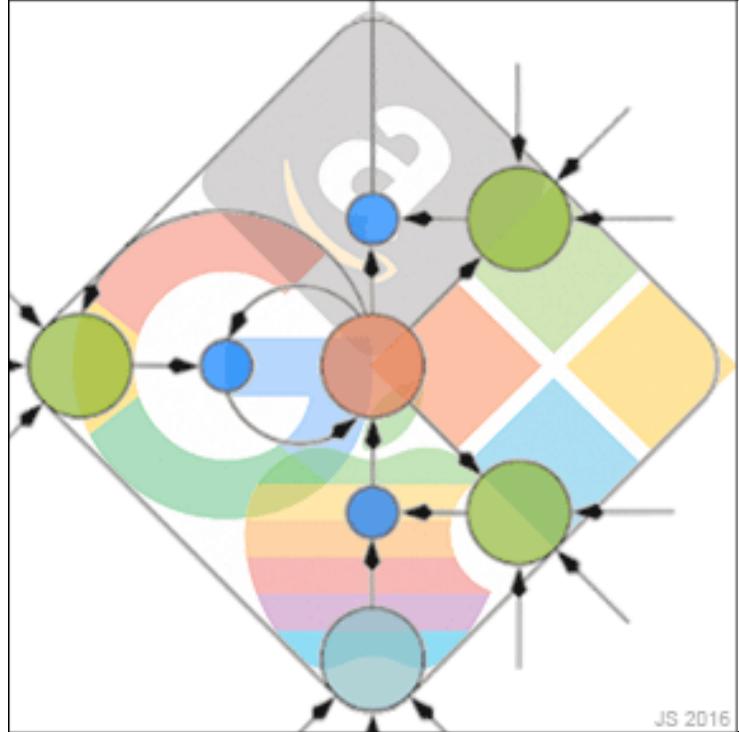
Ex-PhD students (TUM & IDSIA)
[Sepp Hochreiter](#) (PhD 1999), [Felix Gers](#) (PhD 2001, forget gates for recurrent units), [Alex Graves](#) (e.g., CTC, PhD 2008), [Daan Wierstra](#) (PhD 2010), [Justin Bayer](#) (2009, evolving LSTM-like architectures)

Postdocs at IDSIA (2000s)

Fred Cummins, Santiago Fernandez, Faustino Gomez



LSTM cell



JS 2016

Almost 30% of the awesome computational power for inference in all those Google datacenters is now used for LSTM (Jouppi et al, 2017); 5% are used for CNNs discussed later

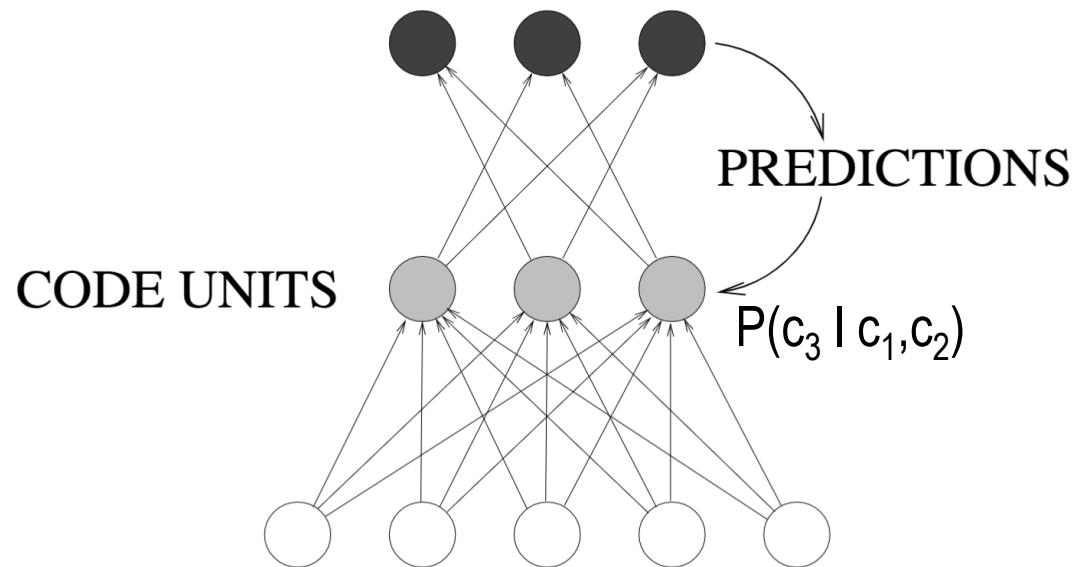


2015: Dramatic improvement of Google's speech recognition through our LSTM & CTC (2006), now on 2 billion Android phones. Similar for Microsoft. 2016: LSTM on almost 1 billion Apple iPhones, e.g., Siri. 2016: Google's greatly improved Google Translate uses LSTM; also Amazon's Echo. 2017: Facebook uses LSTM for over 4 billion translations each day

LSTM / CTC
also used by

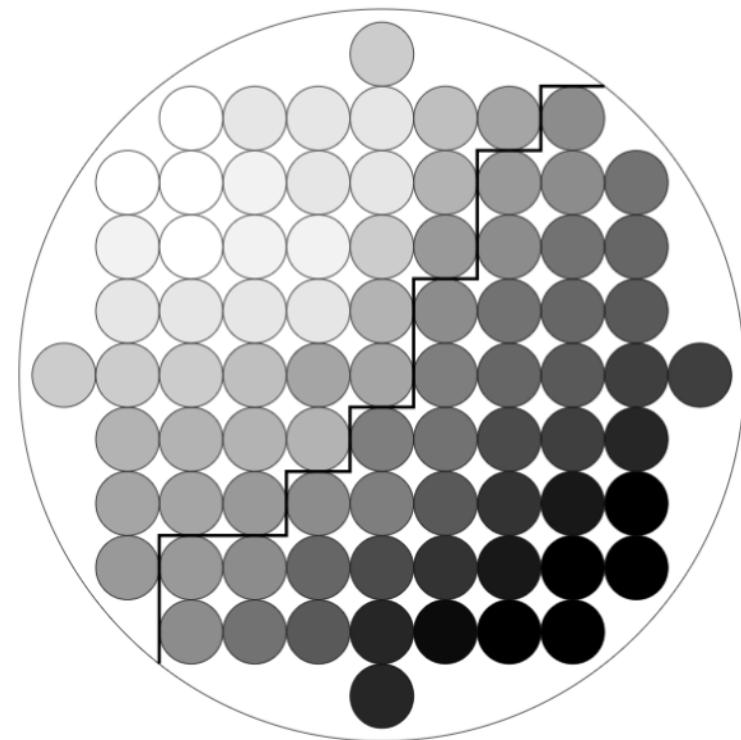



1991: Predictability Minimization (PM): 2 unsupervised nets fight minimax game to model given data distribution

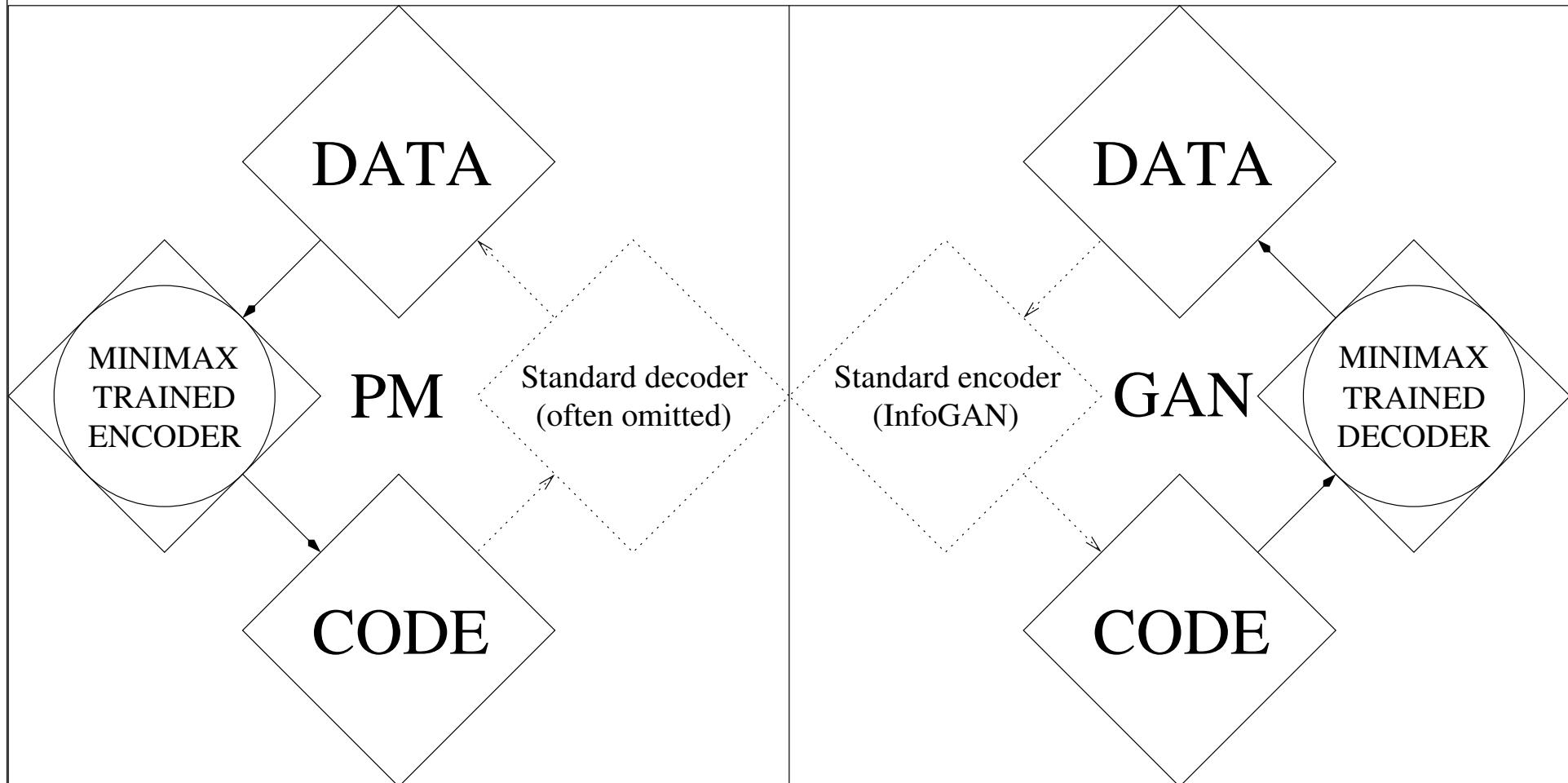


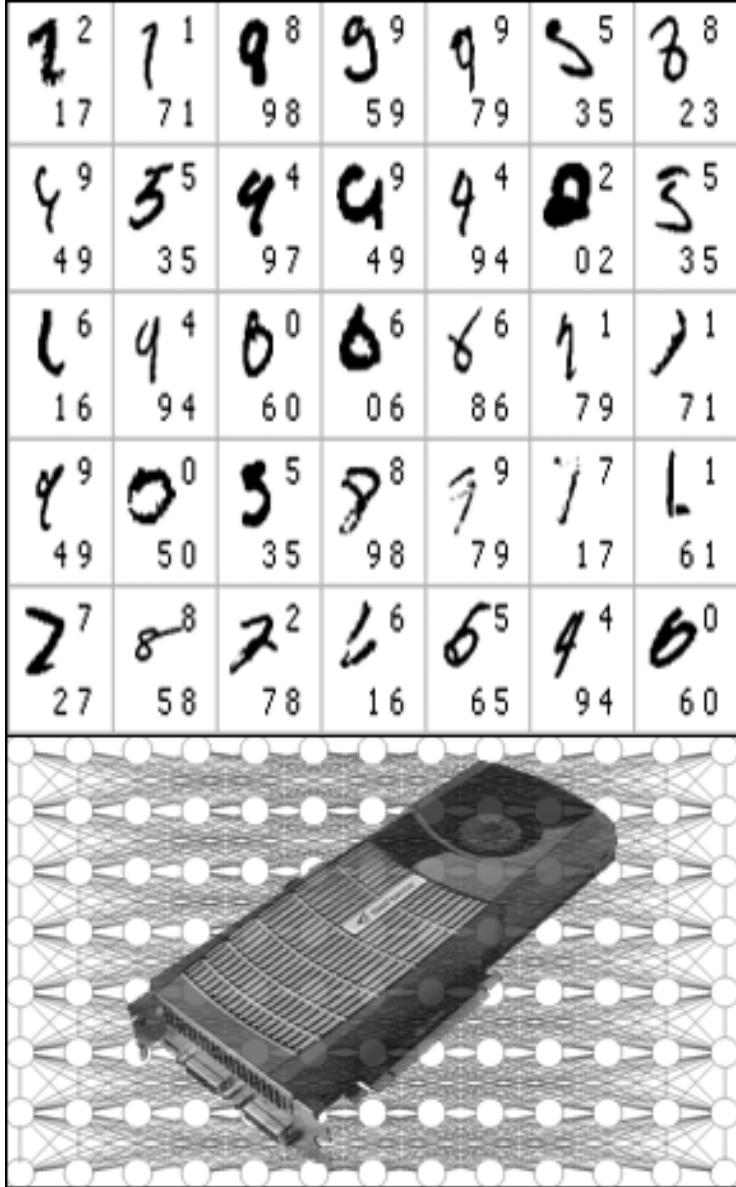
Encoder maximizes objective minimized by predictor. Saddle point = ideal factorial code: $P(\text{pattern}) = P(c_1)P(c_2)\dots P(c_n)$

1996: PM applied to images: learns orientation-sensitive bar detectors, on-center-off-surround detectors, etc



PM v GAN: latent space v original data space





But again, in 2010, pure supervised learning took over, also for feedforward NNs, like for RNNs in the 1990s! 2010: plain backprop for 7 layer MLP, no unsupervised pre-training. MNIST: 60000 digits for training, 10000 for testing; >12m weights; train 200 days on CPU = 5 on GPU; > 10^{15} weight updates, 5b/s, new world record 0.35% ([Neural Comp. 2010, Ciresan, Meier, Gambardella, Schmidhuber](#))

Konrad Zuse 1941
First working
general computer

Every 5 years
10 times cheaper
 $75 \text{ years} \approx 10^{15}$

<http://people.idsia.ch/~juergen/zuse.html>





2011: Traffic Sign Contest, Silicon Valley
Our GPU-CNN was twice better than humans
3 times better than closest artificial competitor
6 times better than best non-neural thing: **FIRST**

SUPERHUMAN VISUAL PATTERN RECOGNITION

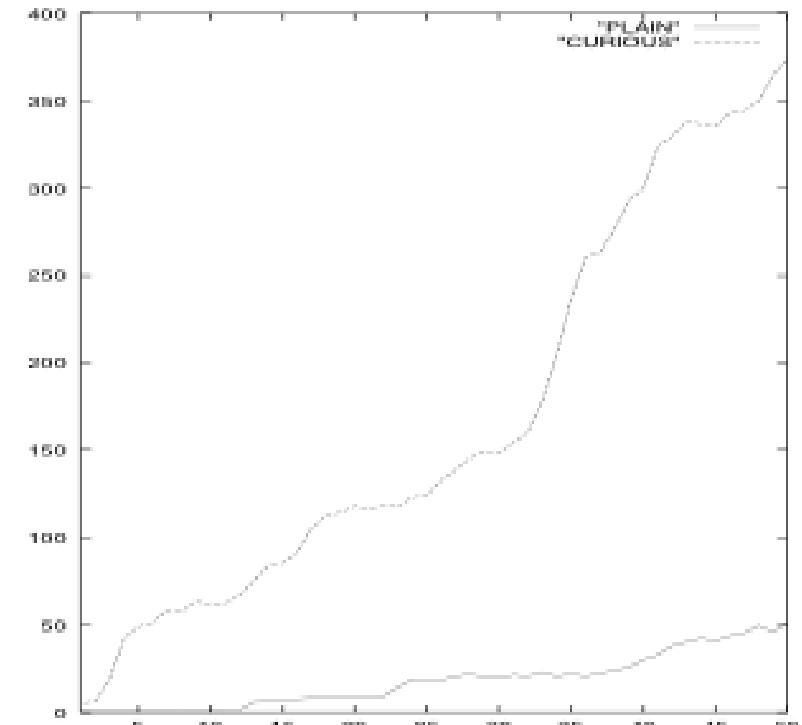
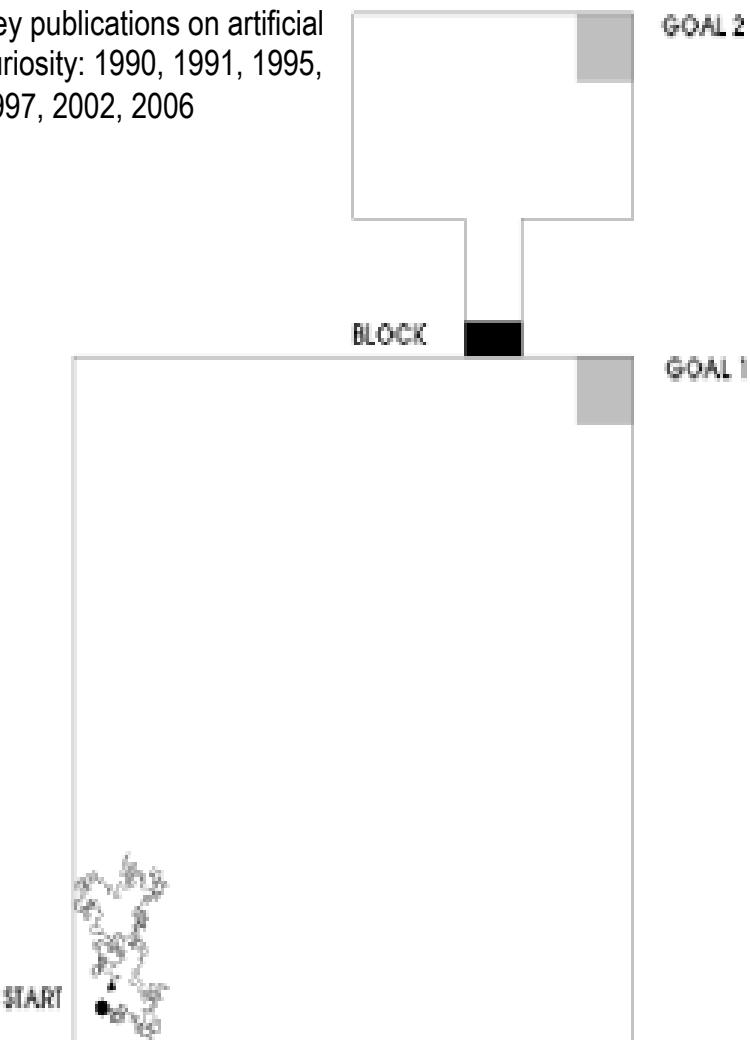
<http://people.idsia.ch/~juergen/superhumanpatternrecognition.html>

JÜRGEN SCHMIDHUBER 2013

1997-2002: Active Unsupervised Minimax for Reinforcement Learning (RL): What's interesting? Exploring the predictable - <http://people.idsia.ch/~juergen/interest.html>

Two reinforcement learning adversaries called "left brain" and "right brain" are intrinsically motivated to outwit or surprise the other by proposing an experiment such that the other agrees on the experimental protocol but disagrees on the predicted outcome, an internal abstraction of complex spatio-temporal events generated through the execution the self-invented experiment. After execution, the surprised loser pays a reward to the winner in a zero sum game. This motivates the two brain system to focus on the "interesting" things, losing interest in boring aspects of the world that are consistently predictable by both brains, as well as seemingly random aspects of the world that are currently still hard to predict by any brain. This type of artificial curiosity can help to speed up the intake of external reward.

Key publications on artificial curiosity: 1990, 1991, 1995, 1997, 2002, 2006



1997-2002: artificial curiosity through active unsupervised minimax accelerates real reward

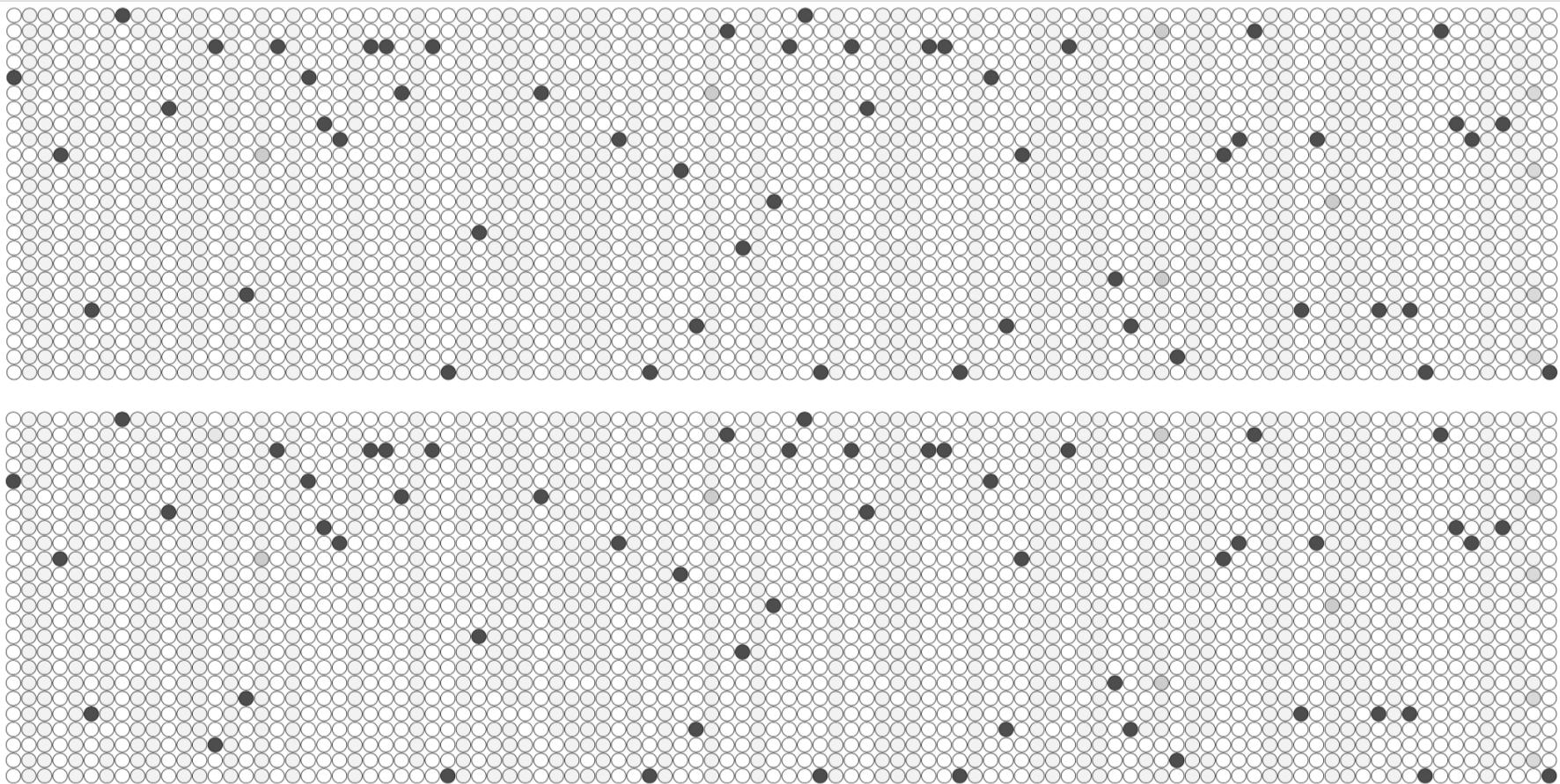


Figure 8: *Experiment 2a: LEFT's (top) and RIGHT's first 100 (of 576) probability distributions after simulation 1. Grey scales indicate probability magnitudes (white = close to 0, black = close to 1). The probability mass of many (but not all) columns is concentrated in a single value. Both brains are almost identical due to SSAandCopy PLAs. Their stacks are quite different though.*

But curiosity can also
kill the cat, and others



J.SCHMIDHUBER

[pm1] J. Schmidhuber. Learning factorial codes by predictability minimization. *Neural Computation*, 4(6):863-879, 1992. Based on TR CU-CS-565-91, Univ. Colorado at Boulder, 1991.

[pm2] J. Schmidhuber, M. Eldracher, B. Foltin. Semilinear predictability minimization produces well-known feature detectors. *Neural Computation*, 8(4):773-786, 1996.

[int1] J. Schmidhuber. What's interesting? TR IDSIA-35-97, IDSIA, July 1997. (Co-evolution of unsupervised RL adversaries in a zero sum game for exploration. See also [int3].)

[int2] J . Schmidhuber. Artificial Curiosity Based on Discovering Novel Algorithmic Predictability Through Coevolution. In P. Angeline, Z. Michalewicz, M. Schoenauer, X. Yao, Z. Zalzala, eds., *Congress on Evolutionary Computation*, p. 1612-1618, IEEE Press, Piscataway, NJ, 1999. Based on [int1].

[int3] J. Schmidhuber. Exploring the Predictable. In Ghosh, S. Tsutsui, eds., *Advances in Evolutionary Computing*, p. 579-612, Springer, 2002. Based on [int1].

More on Predictability Minimization (PM): <http://people.idsia.ch/~juergen/ica.html>

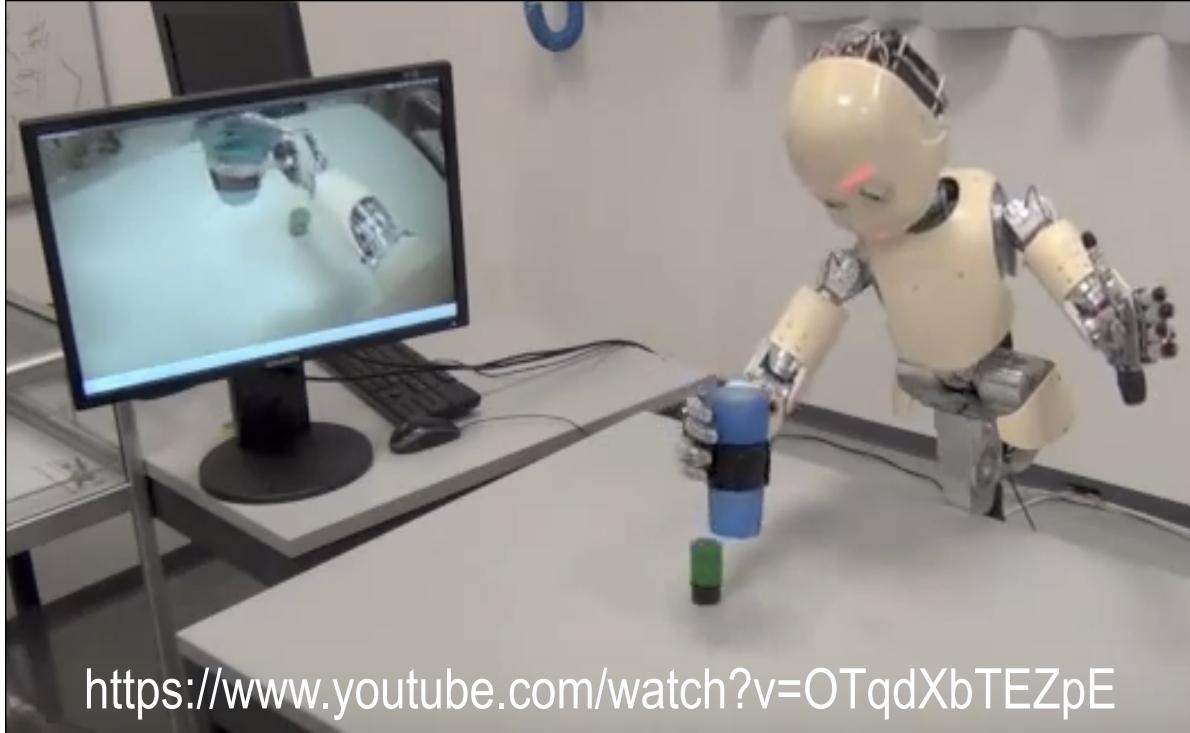
More on artificial curiosity: <http://people.idsia.ch/~juergen/interest.html>

<http://people.idsia.ch/~juergen/creativity.html>

PowerPlay not only solves but also continually invents problems at the borderline between what's known and unknown - training an increasingly general problem solver by continually searching for the simplest still unsolvable problem

POWER PLAY





<https://www.youtube.com/watch?v=OTqdXbTEZpE>

Continual curiosity-driven skill acquisition from high-dimensional video inputs for humanoid robots.
Kompella, Stollenga, Luciw, Schmidhuber. [Artificial Intelligence, 2015](#)

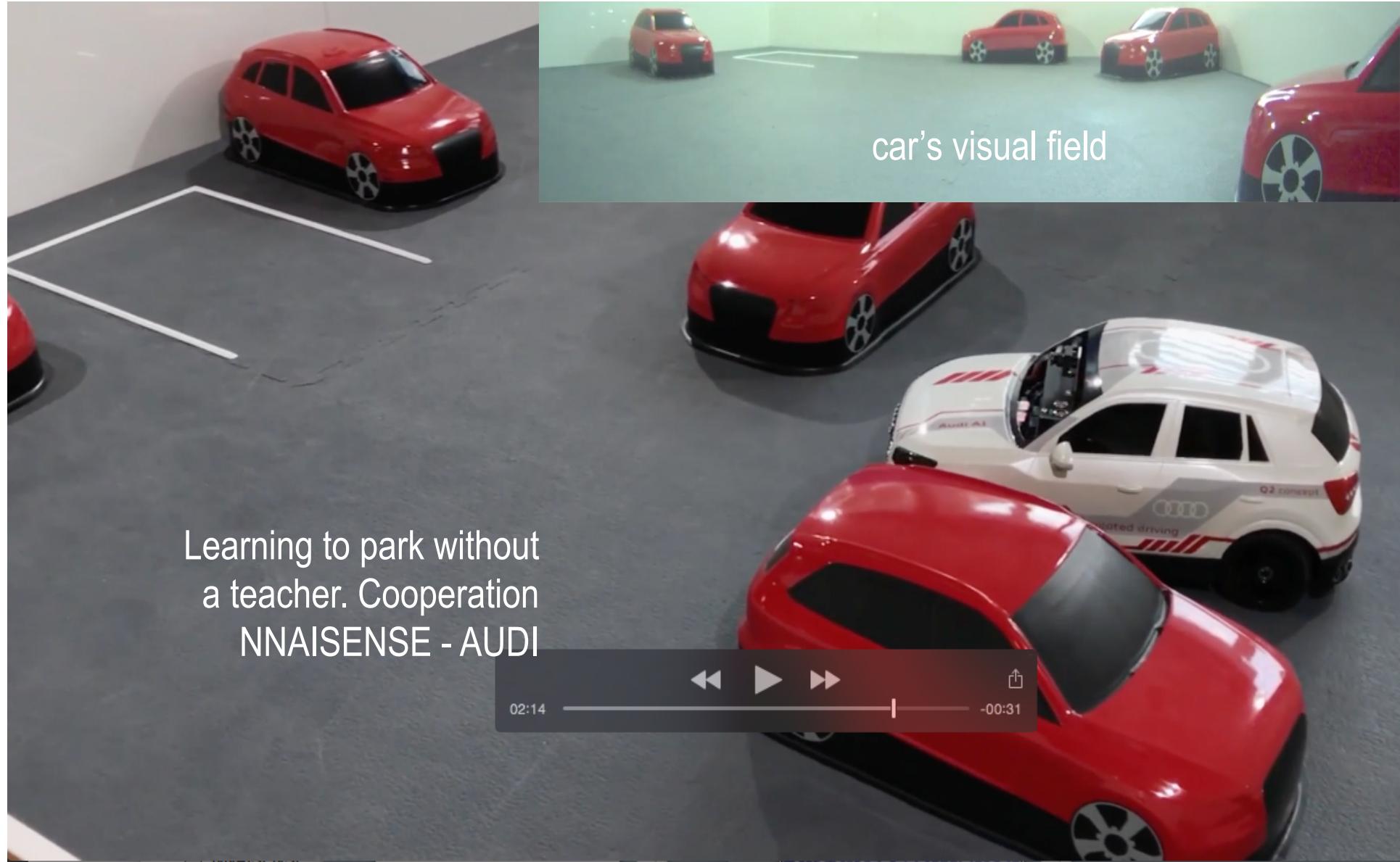




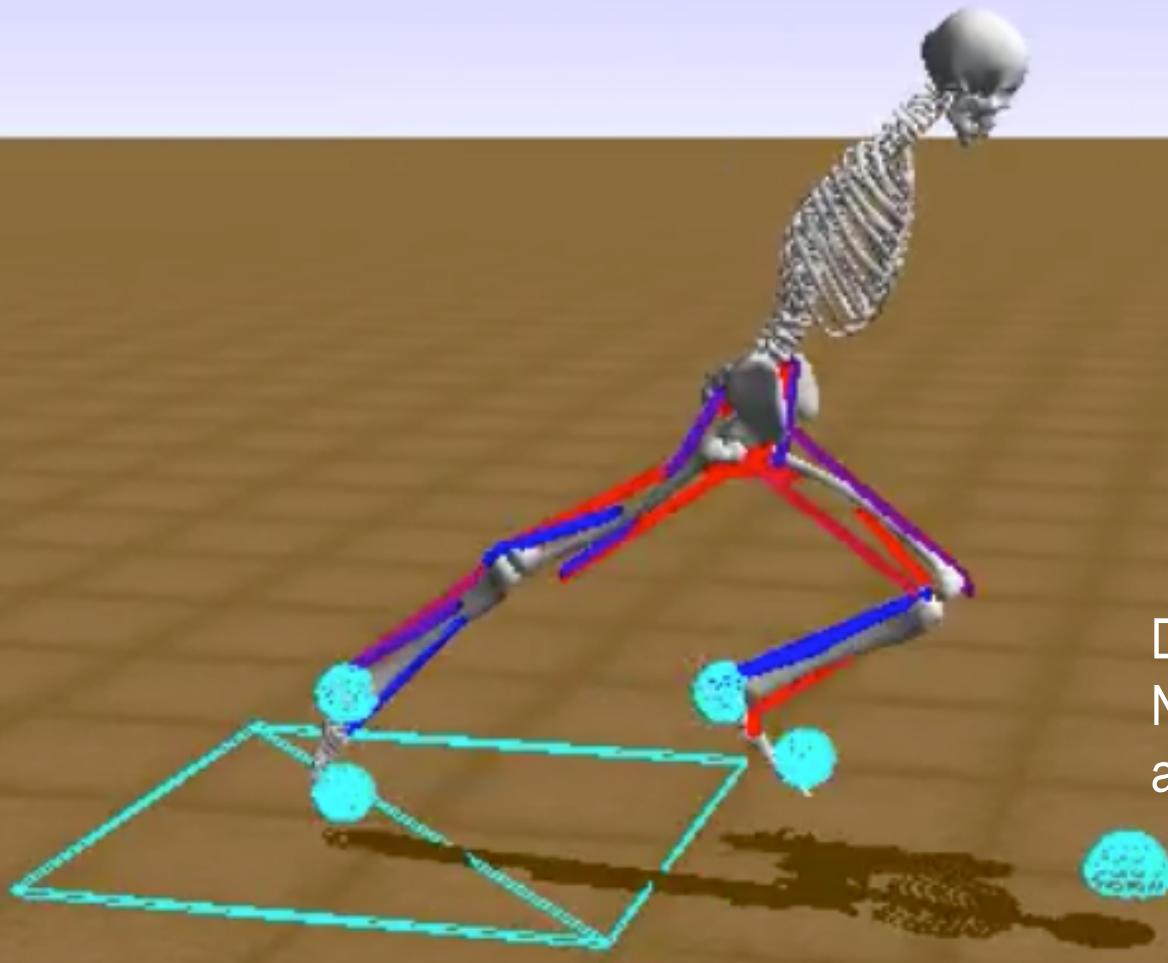
nnaisense

neural networks-based
artificial intelligence

THE DAWN OF AI



car's visual field



Dec 2017: NNAISENSE wins
NIPS “Learning to Run” contest
against over 400 competitors



<http://people.idsia.ch/~juergen/erc2017.html>

www.nnaisense.com