# *Computing goes Light:*
## *Effiziente Datenverarbeitung auf Basis der Photonik*

ADV TRENDS Gamechanger IT

Nov $\frac{30}{}$ 2023

Dr. Bernhard Schrenk

AIT Austrian Institute of Technology

**Computing goes Light:**
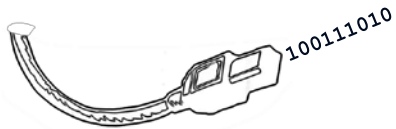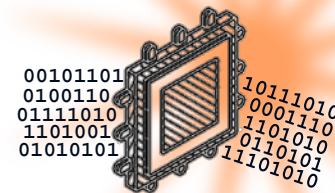**Effiziente Datenverarbeitung auf Basis der Photonik**
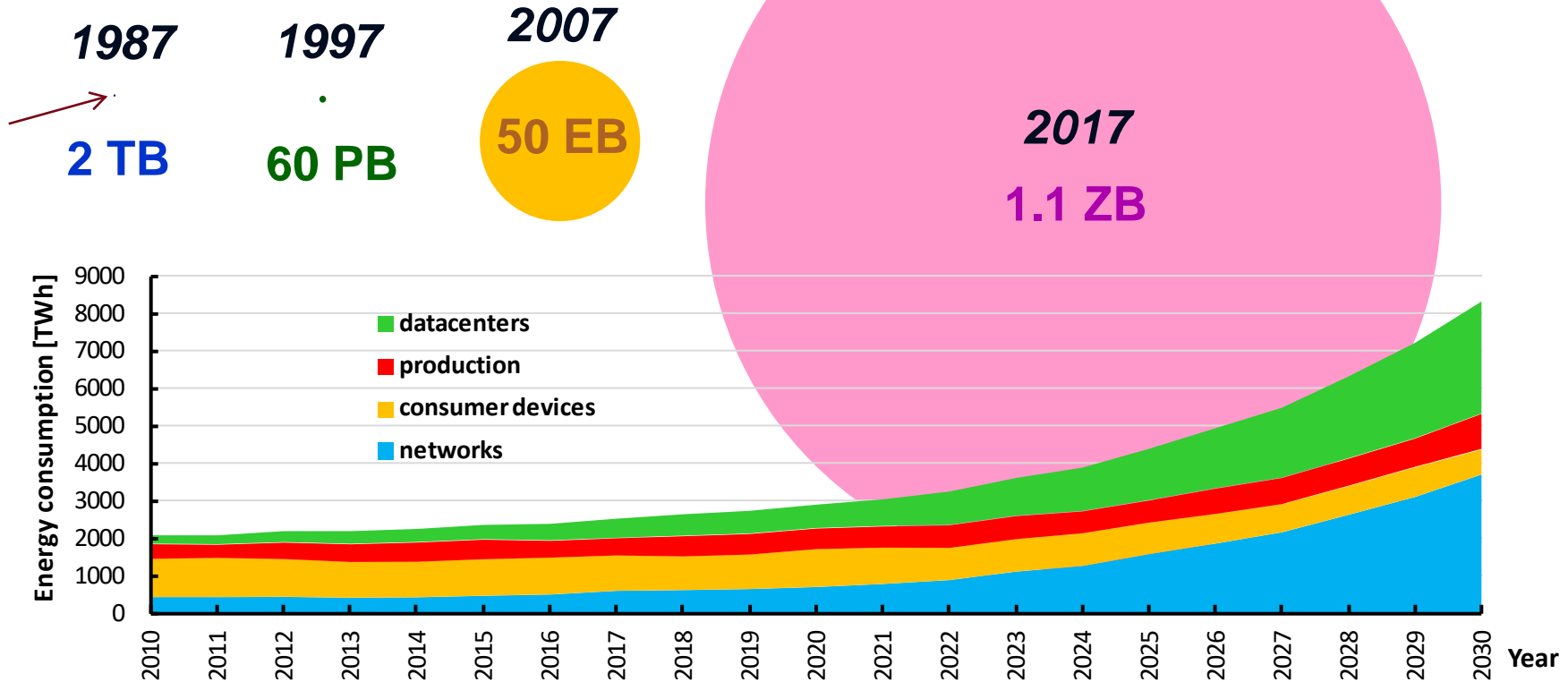
**Communication:**
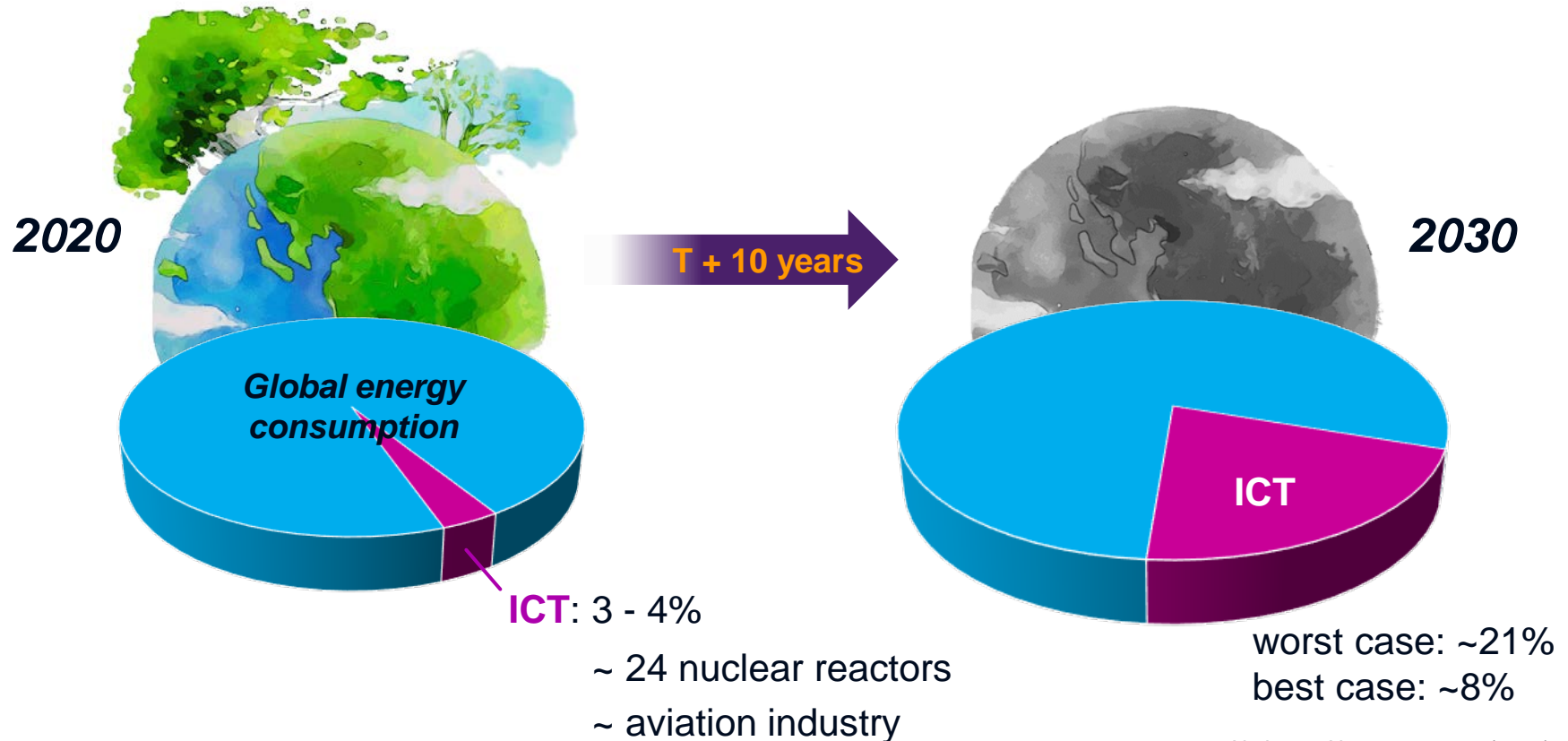**Information in Motion**

**Storage:**
**Information at Rest**

**Processing:**
**Information Transformation**

# The Internet Explosion

*1987*

*1997*

*2007*

**2 TB**

**60 PB**

**50 EB**

*2017*

**1.1 ZB**



- datacenters
- production
- consumer devices
- networks

Energy consumption [TWh]

Year

N. Jones, Nature 561, 163 (2018)

3

# Share of ICT in Global Energy Consumption

*2020*

*2030*

**T + 10 years**

*Global energy consumption*

**ICT**: 3 - 4%

~ 24 nuclear reactors

~ aviation industry

**ICT**

worst case: ~21%
best case: ~8%
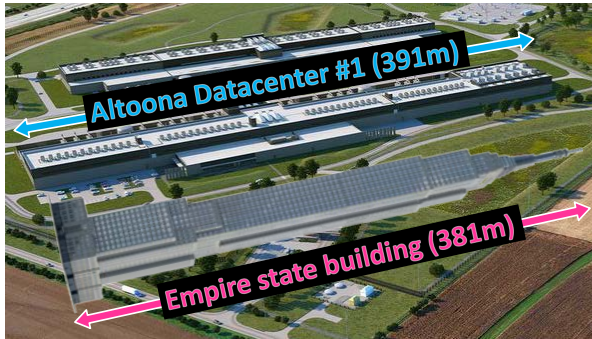
N. Jones, Nature 561, 163 (2018)

# Example: Bitcoin Mining

- 80% happening in China
  - 40% thereof fuelled by coal
- equivalent to 0.6% of world's electricity production
- similar footprint as Italy or Saudi Arabia





**Energy consumption ranking by countries (Twh)**

| Country | Value |
|---|---|
| Canada | 522.2 (No. 7) |
| Brazil | 509.1 (No. 8) |
| South Korea | 507.6 (No. 9) |
| France | 450.8 (No. 10) |
| United Kingdom | 309.2 (No. 11) |
| Bitcoin Industry | 296.6 (No. 12) |
| Saudi Arabia | 296.2 (No. 13) |
| Italy | 293.5 (No. 14) |
| Mexico | 258.7 (No.15) |
| Spain | 239.5 (No.16) |

**Carbon emission ranking by countries (Mton)**

| Country | Value |
|---|---|
| The Netherlands | 166.4 (No. 33) |
| Venezuela | 163.2 (No. 34) |
| Algeria | 147.6 (No. 35) |
| Bitcoin industry | 130.5 (No. 36) |
| Philippines | 119.3 (No. 37) |
| Nigeria | 117.8 (No. 38) |
| Czech Republic | 106.6 (No.39) |
| Qatar | 102.2 (No. 40) |
| Belgium | 98.4 (No. 41) |
| Kuwait | 98.1 (No. 42) |

S. Jiang, Nature Comm. 12, 1938 (2021)

# Processing: Inside the Information Factory

## Cloud Datacenter

Altoona Datacenter #1 (391m)

Empire state building (381m)

~20 000 servers
**~20 MW**

Operating a datacenter: installing a server blade in 2030.

## HPC

Fugaku supercomputer

537 212 TFlop/s
**29.9 MW**

*slow-down*

Demand: Generated data [ZB]

Computing efficiency [MIPS/US$]

Relative transistor density

Year

J. Kendall, Appl. Phys. Rev. 7, 011305 (2020)

## Human

Human brain

~2 000 TFlop/s
**0.000020 MW**

## Compute ops

## el. Power

# Pattern Recognition

Msot plpeoe wlil hvae no peormbls radneig tihs txet, alothguh the oderr of leterts is rndaom (wtih the epeixoctn of the frist and the lsat leettr).

- There is only 1 correct solution and ~ 121 885 070 000 000 000 000 000 possibilities.
- We compute on-the-fly as we read over the text – a fantastic example of pattern recognition.

# AI Hardware

- Multi-layered, deep **neural network**

  - accomodates many **neurons**

    Human brain:   $10^{11}$
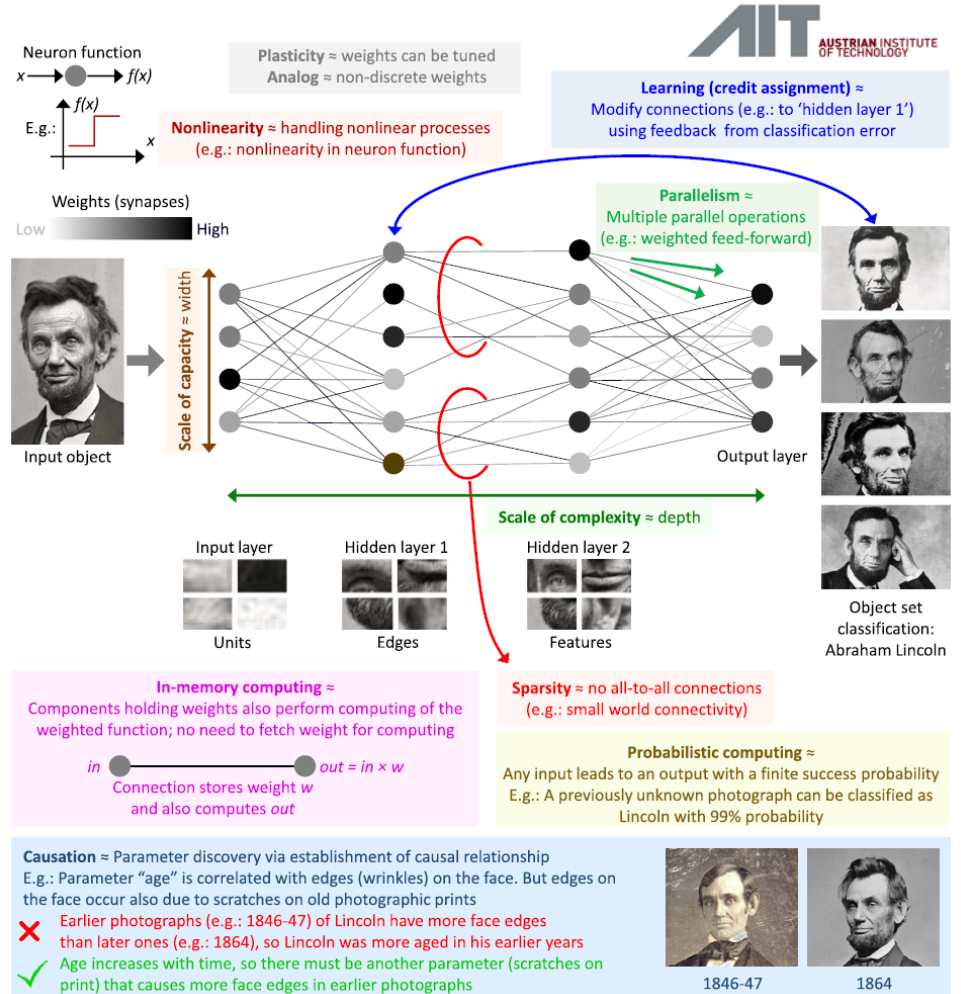    Intel Loihi:       130,000

- Weighted synaptic **interconnect**

  - dense vector-matrix multiplications

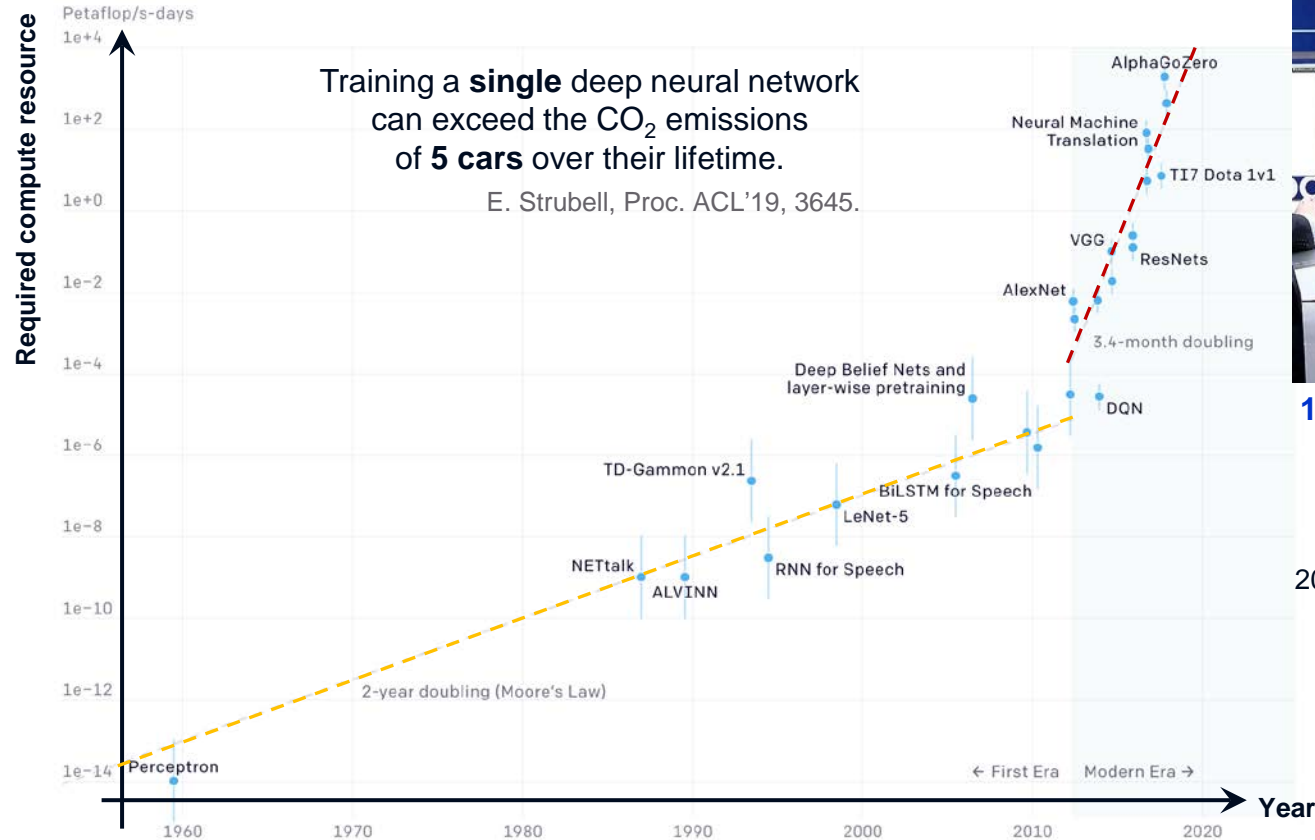  - routing becomes challenging when scaling up the data movement

    Human brain:   $10^4$ inputs/neuron

- Each layer needs to be **trained** …

- … to yield **time-of-flight inference**



**Neuron function**
$x \rightarrow f(x)$

E.g.: $f(x)$

**Plasticity** ≈ weights can be tuned
**Analog** ≈ non-discrete weights

**Nonlinearity** ≈ handling nonlinear processes
(e.g.: nonlinearity in neuron function)

**Learning (credit assignment)** ≈
Modify connections (e.g.: to 'hidden layer 1')
using feedback from classification error

Weights (synapses)
Low        High

**Parallelism** ≈
Multiple parallel operations
(e.g.: weighted feed-forward)

Scale of capacity ≈ width

Input object

Output layer

Scale of complexity ≈ depth

Input layer   Hidden layer 1   Hidden layer 2

Units   Edges   Features

Object set
classification:
Abraham Lincoln

**In-memory computing** ≈
Components holding weights also perform computing of the
weighted function; no need to fetch weight for computing

$in \rightarrow out = in \times w$
Connection stores weight $w$
and also computes $out$

**Sparsity** ≈ no all-to-all connections
(e.g.: small world connectivity)

**Probabilistic computing** ≈
Any input leads to an output with a finite success probability
E.g.: A previously unknown photograph can be classified as
Lincoln with 99% probability

**Causation** ≈ Parameter discovery via establishment of causal relationship
E.g.: Parameter "age" is correlated with edges (wrinkles) on the face. But edges on
the face occur also due to scratches on old photographic prints
✗ Earlier photographs (e.g.: 1846-47) of Lincoln have more face edges
than later ones (e.g.: 1864), so Lincoln was more aged in his earlier years
✓ Age increases with time, so there must be another parameter (scratches on
print) that causes more face edges in earlier photographs

1846-47   1864

AUSTRIAN INSTITUTE
OF TECHNOLOGY

# First:  Training the AI

2016: AlphaGo defeated Lee Sedol

Training a **single** deep neural network can exceed the $CO_2$ emissions of **5 cars** over their lifetime.

E. Strubell, Proc. ACL'19, 3645.

**Required compute resource**

Petaflop/s-days

- 1e+4
- 1e+2
- 1e+0
- 1e-2
- 1e-4
- 1e-6
- 1e-8
- 1e-10
- 1e-12
- 1e-14

AlphaGoZero

Neural Machine Translation

TI7 Dota 1v1

VGG

ResNets

AlexNet

3.4-month doubling

Deep Belief Nets and layer-wise pretraining

DQN

TD-Gammon v2.1

BiLSTM for Speech

LeNet-5

NETtalk

RNN for Speech

ALVINN

2-year doubling (Moore's Law)

Perceptron

← First Era    Modern Era →

**Year**

1960   1970   1980   1990   2000   2010   2020

**1202 CPUs**
**176 GPUs**  vs  **1 human brain**

**1 MW**  vs  **20 W**

2023: Kellin Pelrine defeated Go-playing AI

https://openai.com/blog/ai-and-compute/

# Second: Enjoy Inference at Low Latency

Speech and object recognition

Deep surveillance in real-time

**Accelerators** / co-processors
for vector-matrix multiplication
and deep learning inference,

ultra-fast **control**,

intelligent **signal processing**
(wireless, fiber comms, edge computing)

## Hz – kHz

## kHz – MHz

## GHz

Machine learning
with computers:
AI **software**

Neuromorphic electronics:
**Hardware neural networks**

Challenge: **Interconnect**
(capacitive loading, BW, EMI, routing, leakage / energy)

Source: Presto

Neuromorphic photonics:
**Hardware-based**

Challenge: **scaling (PICs),
all-optical NNs**

10

# PhotonICs

## 1986

**Principal idea:**
**using the silicon manufacturing**
**supply chain to produce photonics**

### The quest for an 'optical silicon'

Even if some researchers would disagree that gallium arsenide and its cousins are the best way to go, most agree on the need for a practical material with all the required optoelectronic properties—what Tanguay has called an "optical silicon." **This ideal material should be versatile, stable, easy to work with, manufacturable, reproducible, and cheap.**

*IEEE Spectrum, 1986*

## 2020  Silicon Photonics

SiPh integrated multi-lane transceivers on 200-mm wafer scale.

TCO:   **< 0.3€ per Gb/s**

0.2€ / mm² in shared fab for 10M chips/year
Saturated 200-mm fab: 50M chips/month
50 datacenters equivalent: 1 Gchips or <2 fab years

W. Bogaerts, 2017

# Si PhotonICs



single-polarisation grating coupler fiber ↔ waveguide

fiber

polarisation-diversity grating coupler fiber ↔ waveguide

intersection of coupled macro-ring filters

thermo-optic heaters

10 µm

broadband splitter

optical delay line

10 µm

waveguide

micro-ring resonator

drop

pass

RF contacts ring modulator

metal contacts

# Si Photon**ICs**

# Electronic vs. Photonic Artificial Neural Network
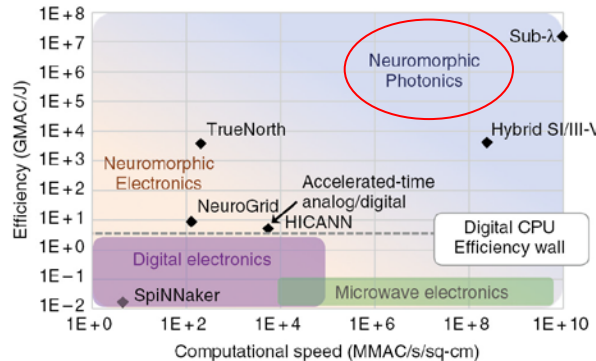
**A neuromorphic processor requires a large number of interconnects!**

**ELECTRONICS**



4096 cores
1 million neurons
256 million synapses
5.4 billion transistors

10 mm

TrueNorth

Neugrid

connections created by
wafer post-processing

HICANN



Power Dissipation (pJ/bit) vs Distance (mm)

- o **Bandwidth-distance trade-off**
- o **Huge energy consumption**
- o **kHz-MHz clock rates**



Photonic neural network

waveguide

input

output

- o **Optical multiplexing**
- o **Low energy consumption**
- o **GHz information rates**

Efficiency (GMAC/J) vs Computational speed (MMAC/s/sq-cm)

Sub-λ ◆

Neuromorphic Photonics

TrueNorth ◆

Hybrid SI/III-V ◆

Neuromorphic Electronics

Accelerated-time analog/digital

NeuroGrid ◆   ◆ HICANN

Digital CPU Efficiency wall

Digital electronics

SpiNNaker ◆           Microwave electronics

**PHOTONICS**



MRR modulator

Photodiode   PUMP   Photodetector (PD)

IN+   Mod.   Heater

IN−   Ge   Photodetector (PD)

Pump   Si WGs   Al wires

OUT   50 µm

Phase shifter

50% coupler

Waveguide   Loss balancing

$M = U\Sigma V$

$V^{(n)}$   $\Sigma^{(n)}$   $U^{(n)}$

Weight SOA matrix

1 mm

Cross-connectivity   De-mux array   Combiners

P. Merolla et. al., Science 345, 668 (2014)

J. Schemmel et al., ISCAS'10 (2010)

B. Benjamin et al., Proc. IEEE 102, 1174 (2014)

Y. Shen et. al., JLT 37, 245 (2019)

T. Ferreira de Lima et al., Nanophot. 6 (2017)

B. J. Shastri et al., Springer (2018)

A. Tait et al., Phys. Rev. Appl. 11 (2019)

Y. Shen et al., Nat. Phot. 11 (2017)

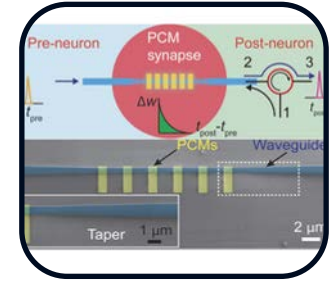B. Shi et. al., JSTQE 26 (2020)

# The Neuron Goes "Light"



$$\sum_k \nu^{(k)} w_k$$

$$w_k \in [-1, +1]$$

**analog linear**

*excitatory*

*inputs*

dendritic tree

$w^+$ soma    axon

$\Sigma$    $\theta$

output

$w^-$

*non-linear*

*inhibitory*

Optical interference unit

SU(4) core    DMMC

Y. Shen et al., Nat. Phot. 11 (2017)

Weight SOA Matrix

Cross connectivity    De-mux array    Combiners

B. Shi et al., JSTQE 26 (2020)

MRR weightbank    Balanced PD
WDM input    RF output
WDM Weighted Addition
IN    THRU
DROP

A. Tait et al., JSTQE 22 (2016)

Pre-neuron    PCM synapse    Post-neuron
Taper 1 μm    2 μm

A. Tait et al., JSTQE 22 (2016)

partial reflection
SOA
EAM    highly reflective coating

M. Stephanie et al., JLT 41 (2023)

$o$

$i$

$\theta$

15

# The Neuron Goes "Light"



$$\sum_k \nu^{(k)} w_k$$

$$w_k \in [-1, +1]$$

B. Shastri et al., Sci. Rep. **6** (2016)

G. Li et al., Nanophotonics (2022)

B. Schrenk, ECOC'19, We.P27 (2019)

*linear*

*excitatory*

$w^+$   soma

**inputs**

dendritic tree

axon

$\theta$

*output*

$w^-$

*inhibitory*

**analog non-linear**

R. Amin et al., APL Materials **7** (2019)

M. Stephanie et al., Proc. SUM, MF4.4 (2023)

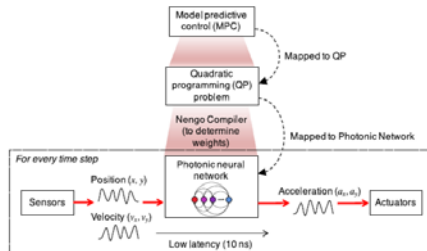# Accuracy at Speed



Iris Classification Problem
150 flower samples

Iris Versicolor    Iris Setosa    Iris Virginica

Input    Layer 1    Layer 2    Layer 3   Output

$x_0$   $x_1$   $x_2$   $x_3$

$A$   $B$   $w_+$   $w_-$   $w_-$   $w_+$   $D$   $E$   $z_0$   $z_1$   $z_2$

Neuron
$\Sigma$
$w_0$   $w_1$   $w_2$   $w_3$   $b$
$f(\Sigma)$
$(+)$   $(-)$
ReLU

NN accuracy

93%   92%   91%   92%

dig.   opt.   dig.   opt.
**weighted sum**    **ReLU**

***Accuracy***

**digital NN: 93%**

**optical NN: 91-92%
but one flower every ns!**

M. Stephanie et al., Proc. SUM, MF4.4 (2023)
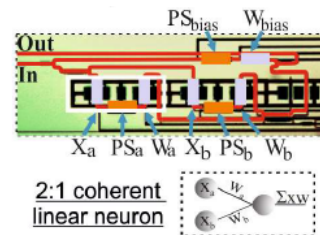
# What about "real" Applications?

**Predictive Control**

❖ for object at flight

❖ 24 neurons
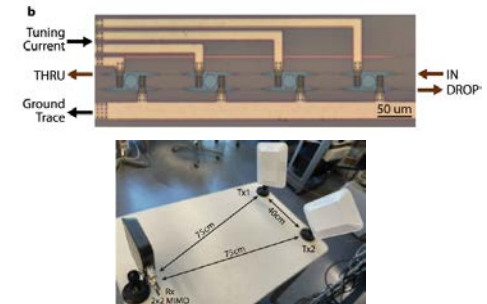
❖ convergence time of 10 ns



**Distributed Denial of Service (DDoS) Attack Identification**

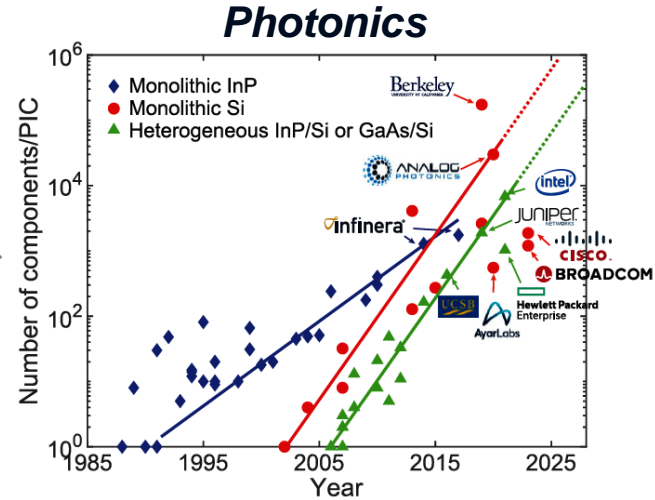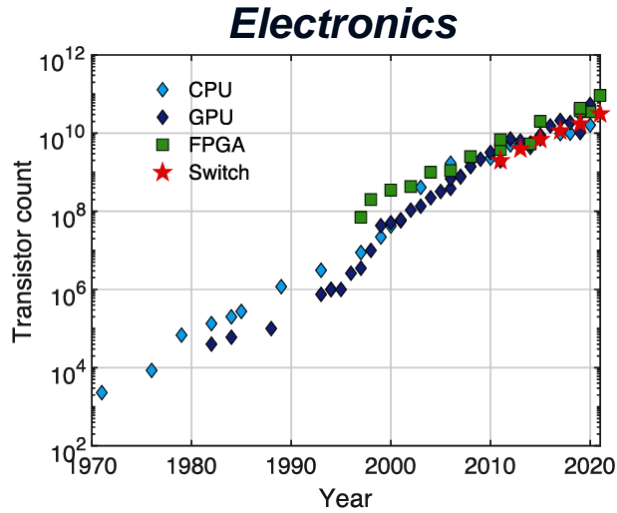❖ Using silicon photonic processor
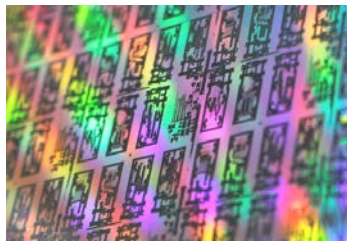
❖ 50 GHz signal rate



**Blind Source Separation**

❖ Separating an unknown mixture of unknown independent signals

❖ Using microring weight bank
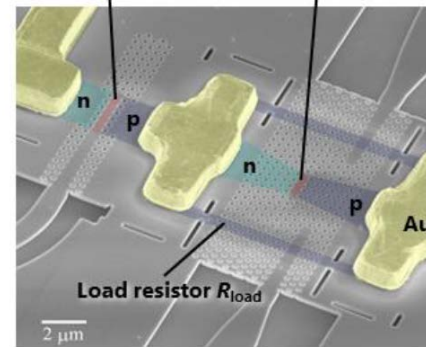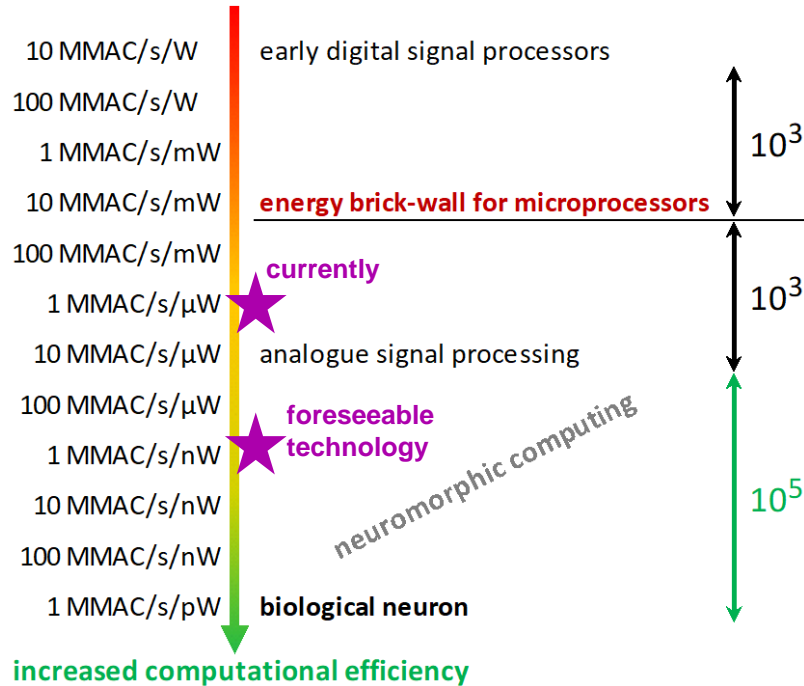
❖ achieving a processing bandwidth of up to 19.2 GHz



T. Ferreira de Lima et. al., JLT 37 (2019)     A. Tsakyridis et al., Proc. OFC (2023)     W. Zhang et. al., Nature Comm. 14 (2023)     18

# Is there a Moore's Law for PICs?



N. Margalit, Appl. Phys. Lett. 118, 220501 (2021)

# …and last, what about Energy Efficiency?



| | |
|---|---|
| 10 MMAC/s/W | early digital signal processors |
| 100 MMAC/s/W | |
| 1 MMAC/s/mW | $10^3$ |
| 10 MMAC/s/mW | **energy brick-wall for microprocessors** |
| 100 MMAC/s/mW | **currently** |
| 1 MMAC/s/µW | $10^3$ |
| 10 MMAC/s/µW | analogue signal processing |
| 100 MMAC/s/µW | **foreseeable technology** |
| 1 MMAC/s/nW | |
| 10 MMAC/s/nW | neuromorphic computing |
| 100 MMAC/s/nW | $10^5$ |
| 1 MMAC/s/pW | **biological neuron** |

**increased computational efficiency**

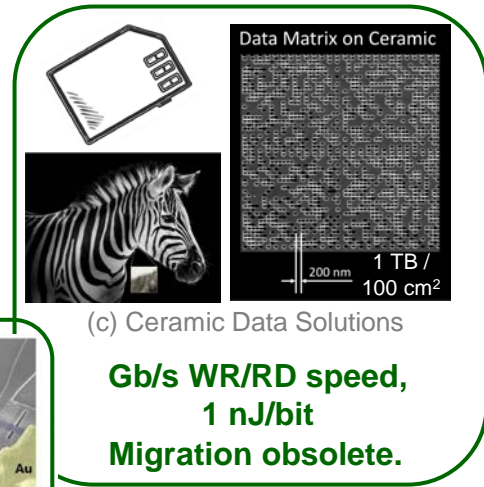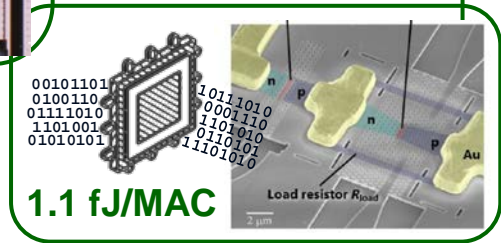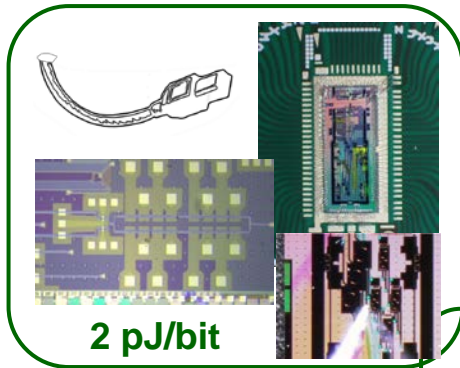K. Nozaki, Nat. Photon. **13**, 454 (2019)

foreseeable:

1.1 fJ/MAC

Lao Tzu

# Take Awayyy

**There are promising solutions to keep the raising ICT energy footprint in check.**

**At the same time, we obtain better performance than the state-of-play.**



**ICT**

**2 pJ/bit**

**1.1 fJ/MAC**

(c) Ceramic Data Solutions

Data Matrix on Ceramic

1 TB / 100 cm²

200 nm

**Gb/s WR/RD speed, 1 nJ/bit Migration obsolete.**

# *Another Decathlon for Another Decade!*

green

bandwidth

flexible

cost-effective

scalable

**Bernhard Schrenk**

bernhard.schrenk@ait.ac.at

Center for Digital Safety & Security
AIT Austrian Institute of Technology

# Communication: Information in Motion

- between people
- between machines

$$\text{Energy / bit} = \frac{\text{Power consumption}}{\text{Data rate}}$$

## Wireless

1 Gb/s

$P_{TX}$ = 100 mW
PAE ~40% (back-off)
OFDM DSP

## Copper

Eth switch

Retimer on backplane

## Optical

Opto-el. transceivers (Si, InP)

**1000 pJ/bit** → LiFi? (10-100 Gb/s)

**18 pJ/bit** → **2 pJ/bit** Active optical cables, PON, etc

# Storage: Information at Rest

- short-term caching
- long-term storage

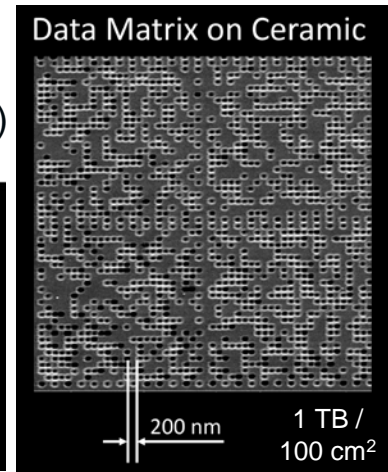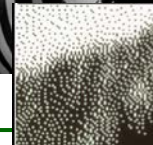$$\frac{Energy}{Preserved\ bit} \quad + \quad Hardware\ Resources$$

## HDD / SSD



## Optical memory

Cold storage: archival (30+ years)



Data Matrix on Ceramic

200 nm

1 TB / 100 cm$^2$

(c) Ceramic Data Solutions

**30-60 nJ/bit/year
migrate every 2-3 years**

**write / read at Gb/s speed and 1 nJ/bit
No migration required.**