# Bit-Tactical: Sparse/Dense Deep Learning Inference Accelerator(*)

Alberto Delmas, Patrick Judd, Mostafa Mahmoud, Milos Nikolic, Zissis Poulos, Sayeh Sharify, Kevin Siu, Dylan Stuart
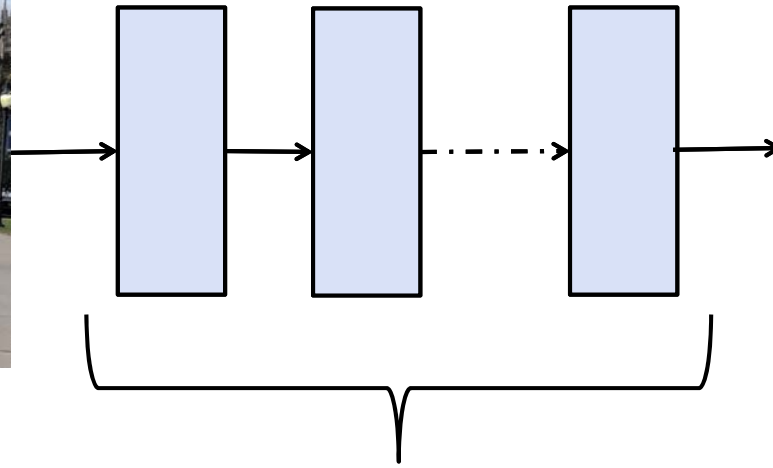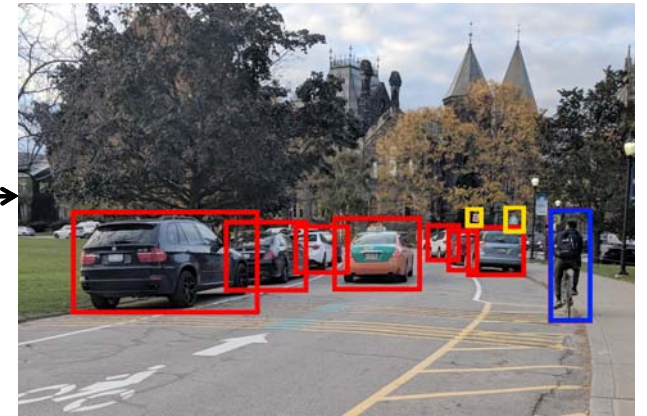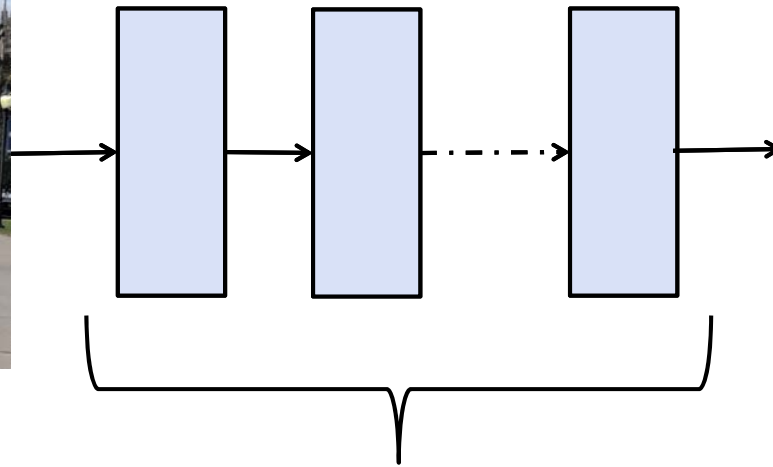
Andreas Moshovos

Input image
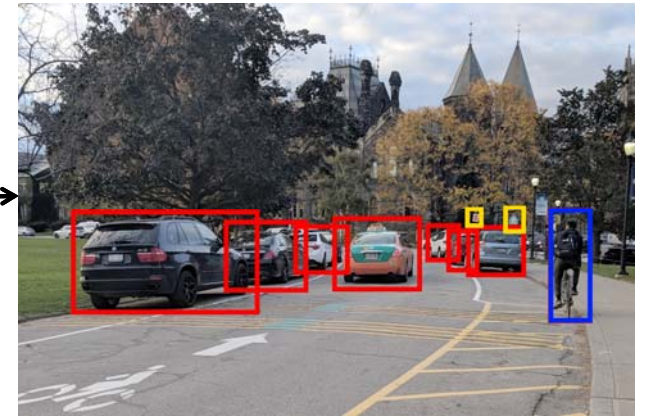
layers
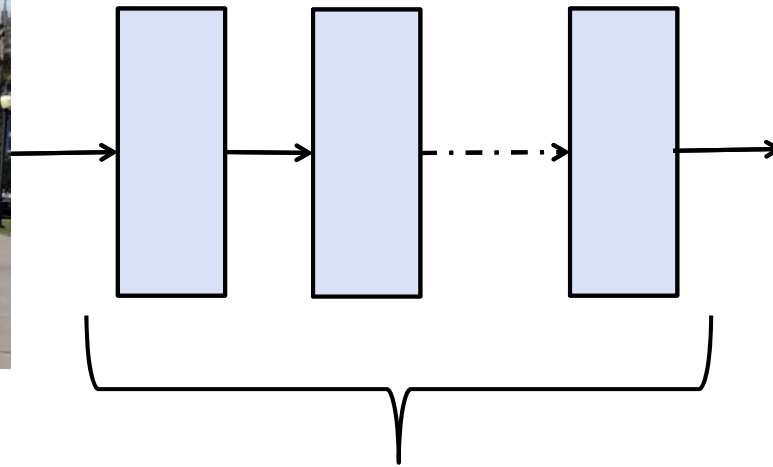10s-100s

Annotated image

Input image

layers
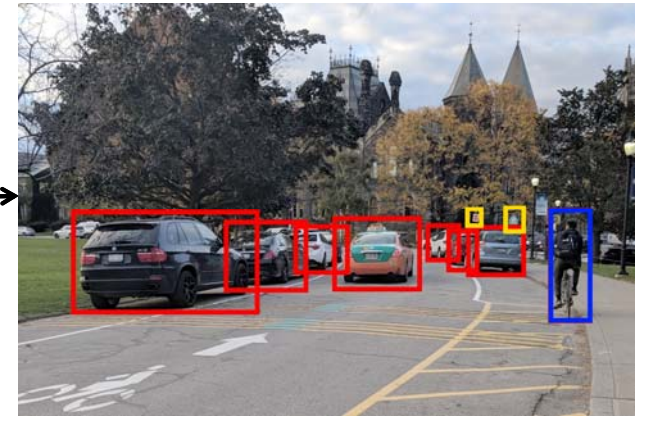10s-100s

Annotated image

Faster
+
More Energy Efficient

Input image

layers
10s-100s

Annotated image

Fewer Operations
Fewer Data Transfers

# Fewer Operations
# Fewer Data Transfers

# Fewer Operations
# Fewer Data Transfers

$$Out\ +=\ A_0 \times W_0$$

$$Out\ +=\ A_1 \times W_1$$

$$Out\ +=\ A_2 \times W_2$$

$$Out\ +=\ A_3 \times W_3$$

$$Out\ +=\ A_4 \times W_4$$

:

# Eliminate Ineffectual Operations

$Out \mathrel{+}= A_0 \times W_0$

$Out \mathrel{+}= A_1 \times W_1$

$Out \mathrel{+}= A_2 \times W_2$

$Out \mathrel{+}= A_3 \times W_3$

$Out \mathrel{+}= A_4 \times W_4$

:

# Eliminate Ineffectual Operations

Out += $A_0$ x $W_0$

Out += 0 x $W_1$

Out += $A_2$ x $W_2$

Out += $A_3$ x 0

Out += $A_4$ x $W_4$

# Takeaway #1

**0.5x – 0.2x**  <span style="color:red">Off-Chip Transfers</span>
<span style="color:red">On-Chip Storage</span>

*DPRed: Making Typical Activation and Weight Values Matter In Deep Learning Computing*, Delmas et al.**,**
**https://arxiv.org/abs/1804.06732**

# Takeaway #2

## Out += A x 0

# Takeaway #2

**Out += A x 0**

Aim to get **most** not **all**

Simple Design
Software Scheduler

# Takeaway #3

## Out += 0 x W

Don't!

Only ~50% A==0

# Takeaway #3

**Out += 0100001 x W**

Go for Bit Sparsity

> 90%

# Takeaway #4

## vs. SCNN

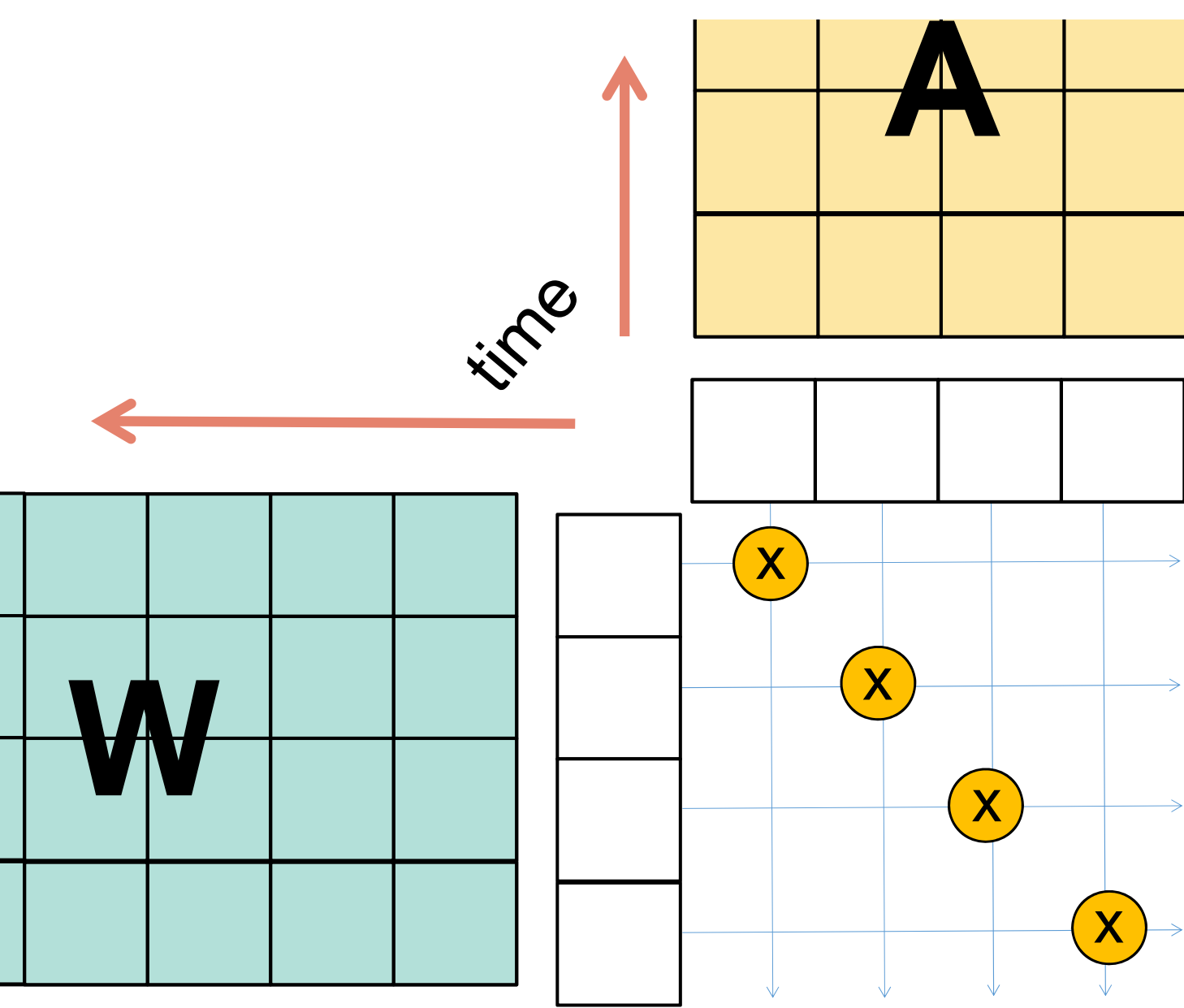**Speed**
5.3x

**Energy**
0.85x

Do as you are told?

$$Out \mathrel{+}= A_0 \times W_0$$

$$Out \mathrel{+}= A_1 \times W_1$$

$$Out \mathrel{+}= A_2 \times W_2$$

$$Out \mathrel{+}= A_3 \times W_3$$

$$Out \mathrel{+}= A_4 \times W_4$$

$$\vdots$$

A

W

time

x
x
x
x

time

time

2 steps in time

Another 2 steps in space

"unrestricted" motion

A: Dynamic Sparsity Pattern

# Cambricon-X



time
"any"

space
"any"

# Nvidia's SCNN



time
"any"

space
"any"

# Bit-Tactical



Weights

time
1-2 steps

space
2-5 places

Do as you are told?

Out += $A_0$ x $W_0$

Out += $A_2$ x $W_2$

Out += 0 x $W_3$

Out += $A_4$ x $W_4$
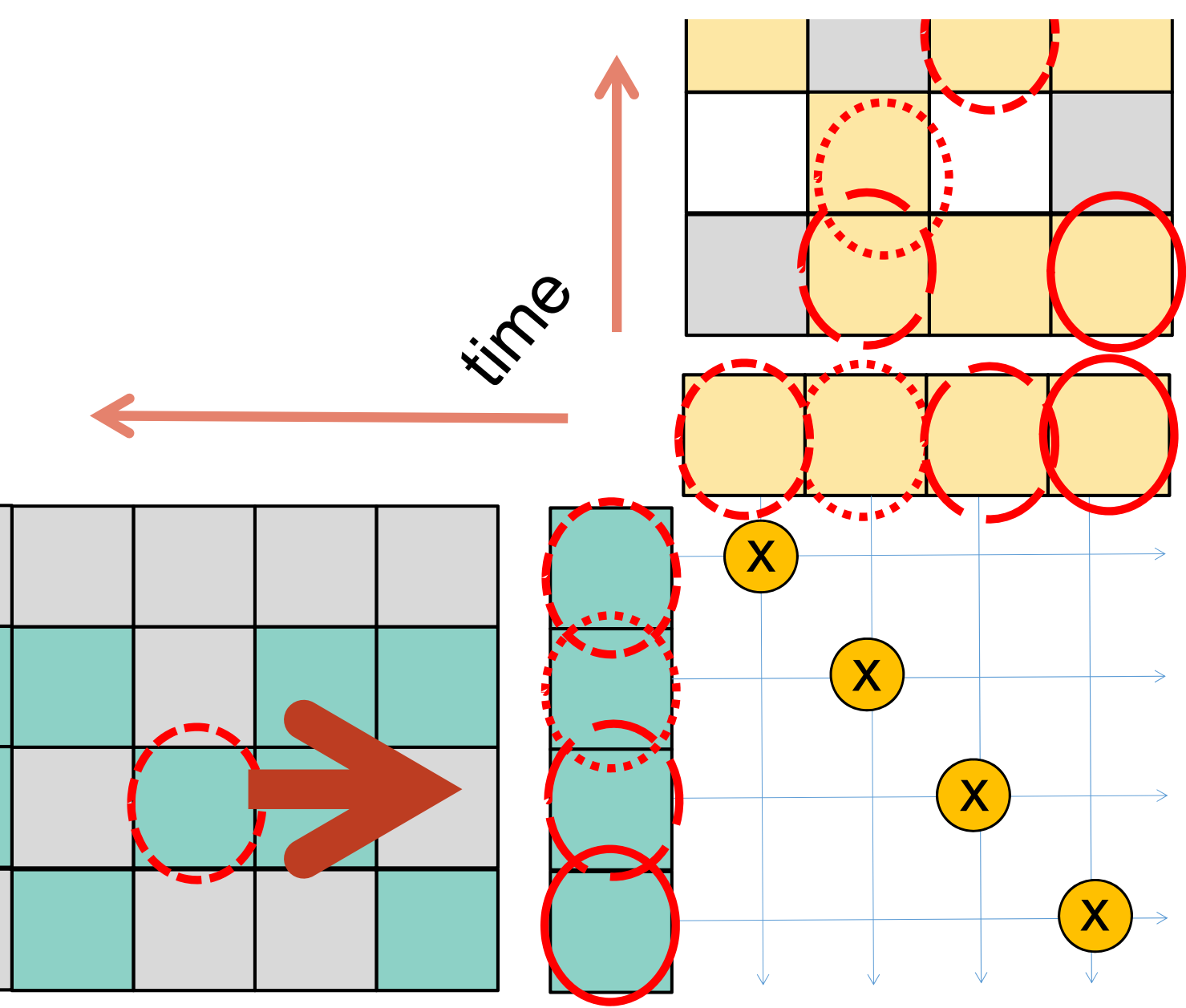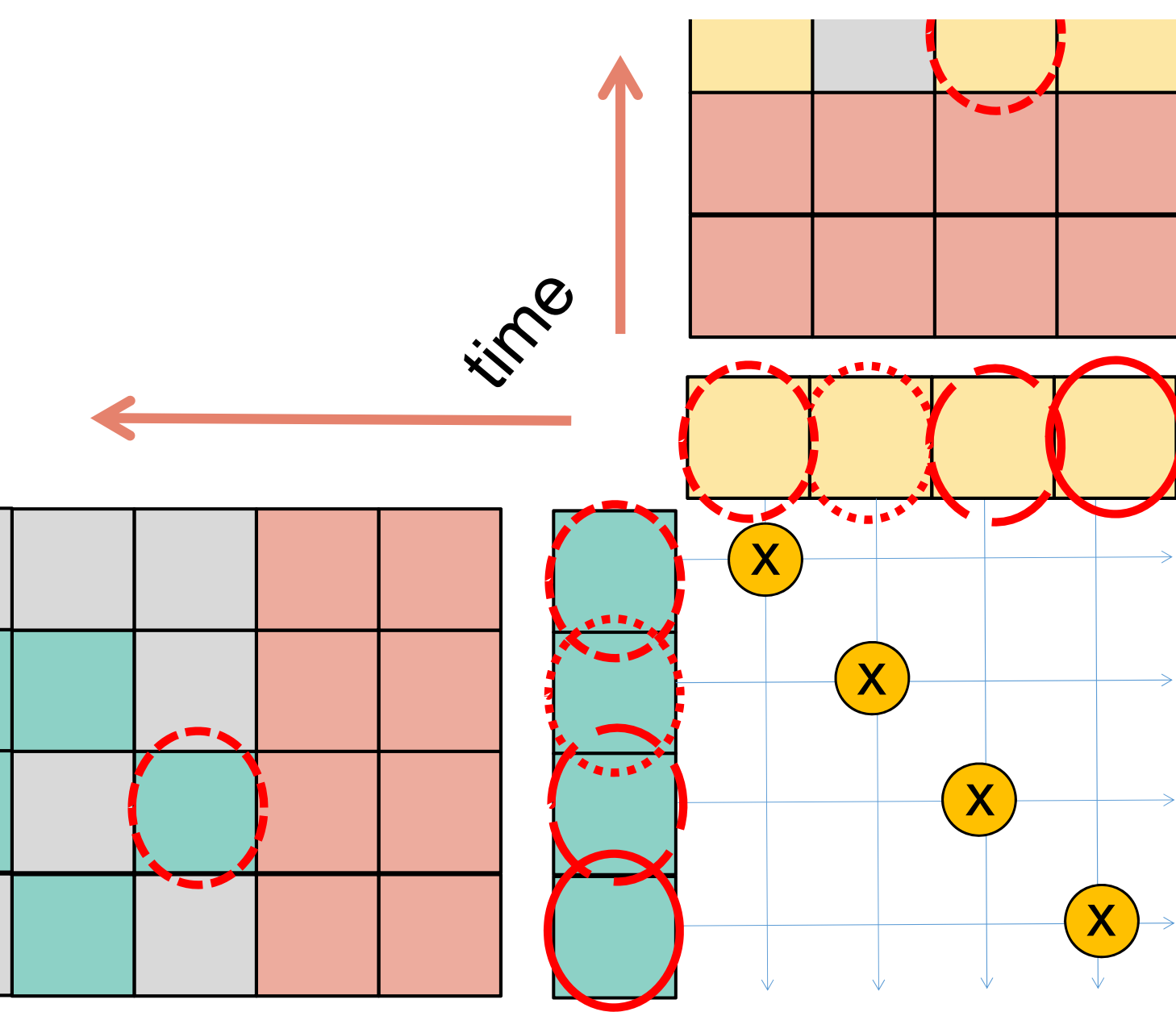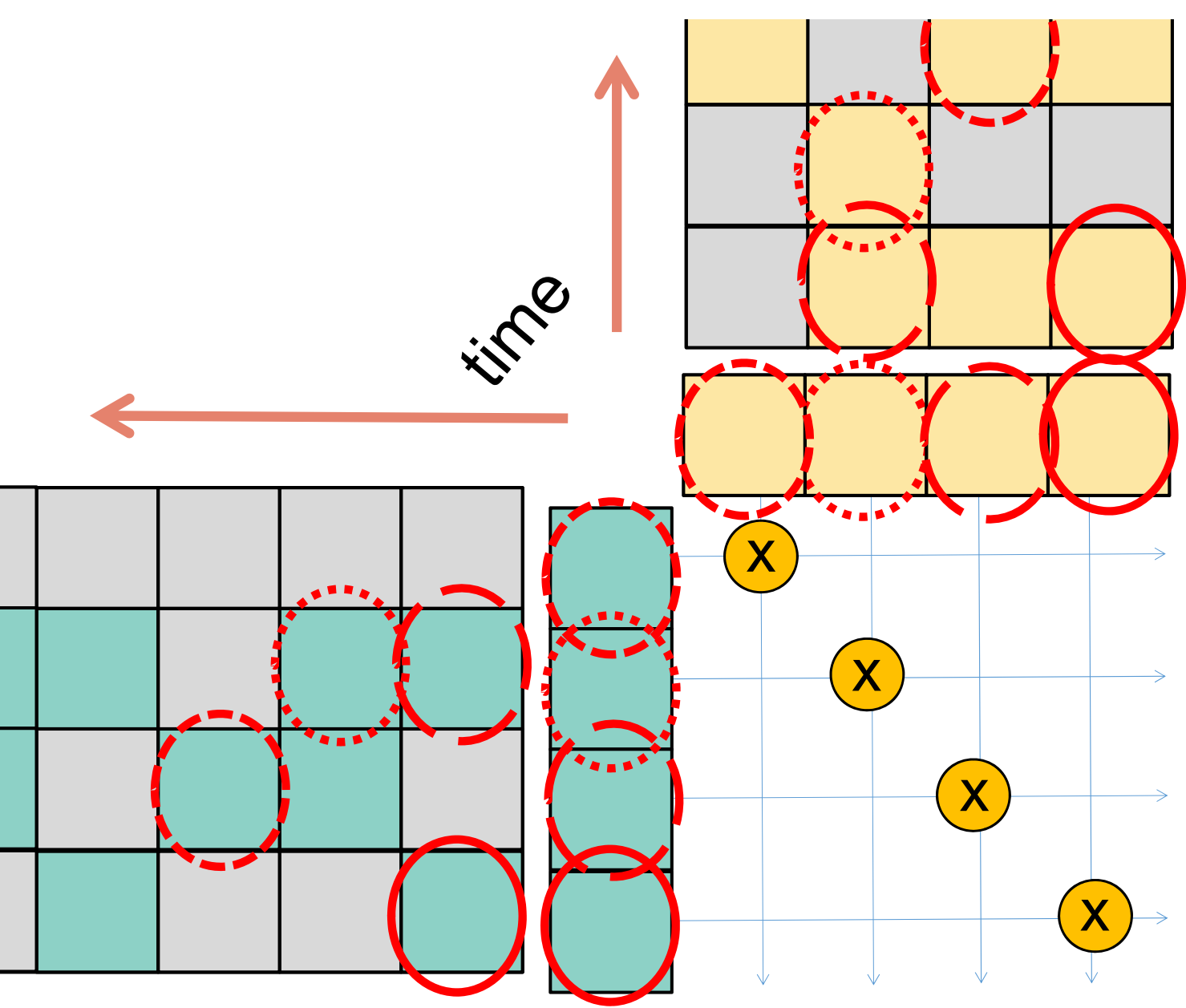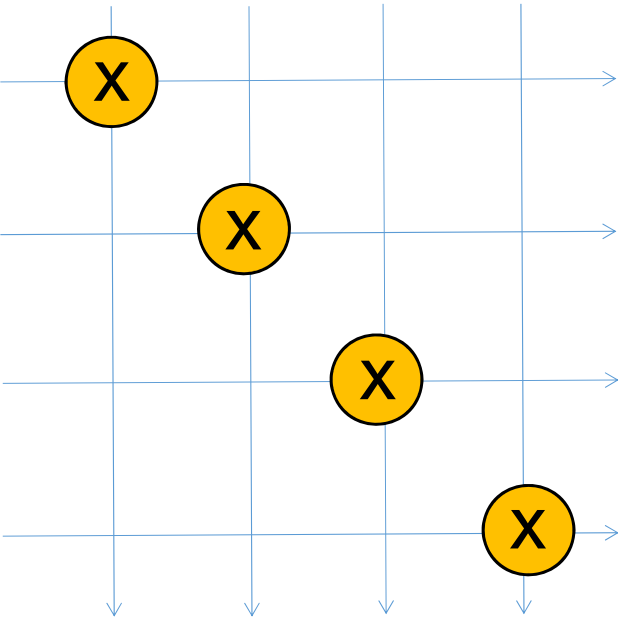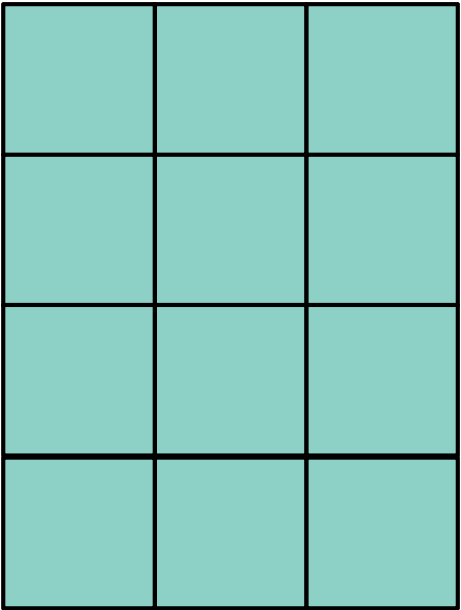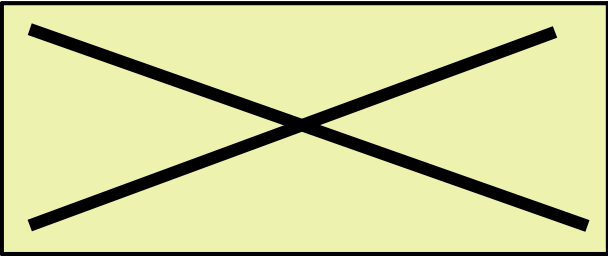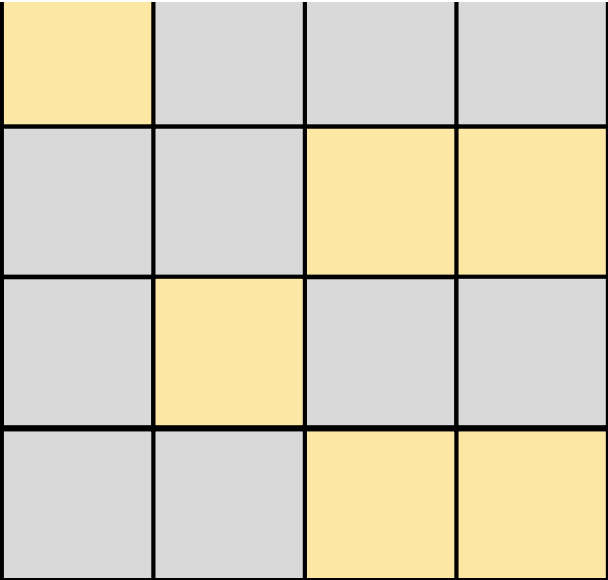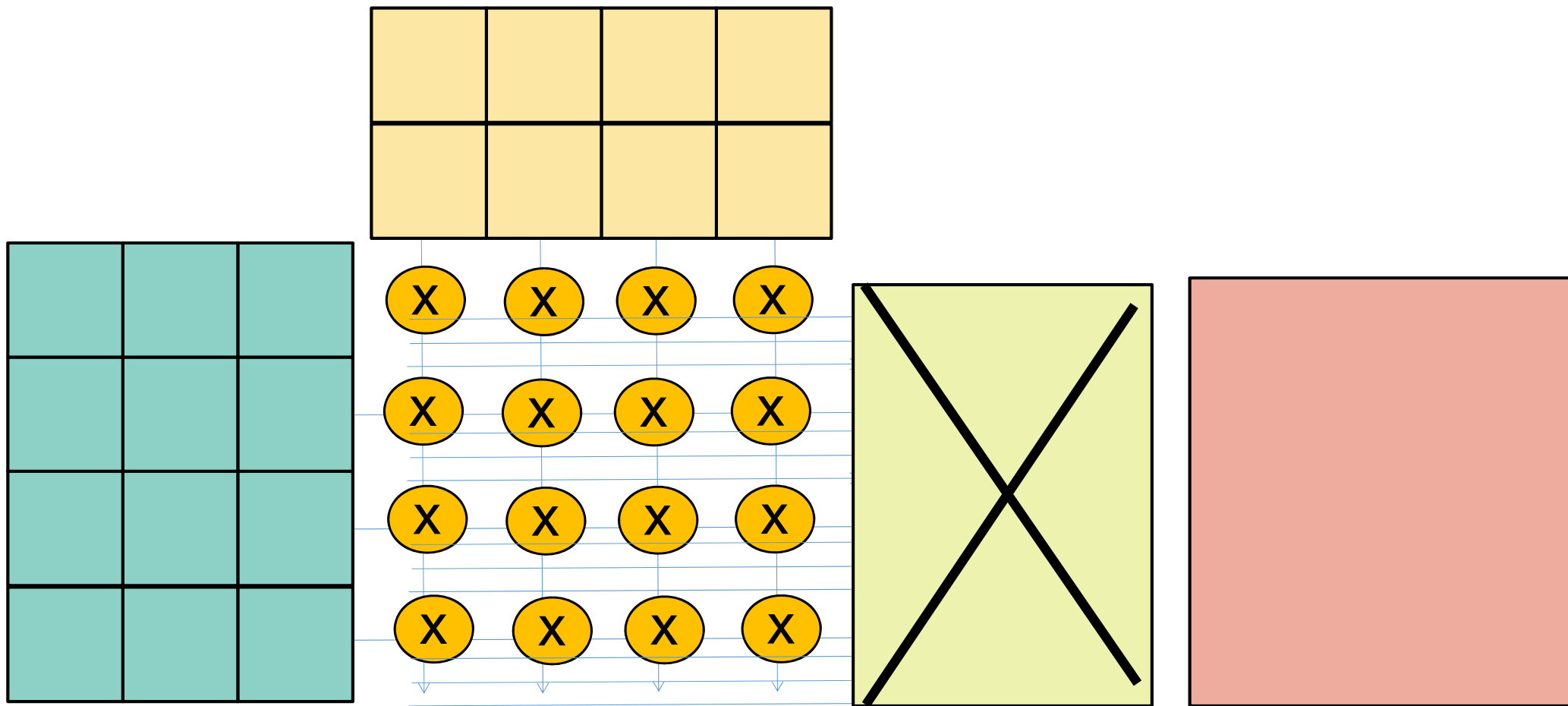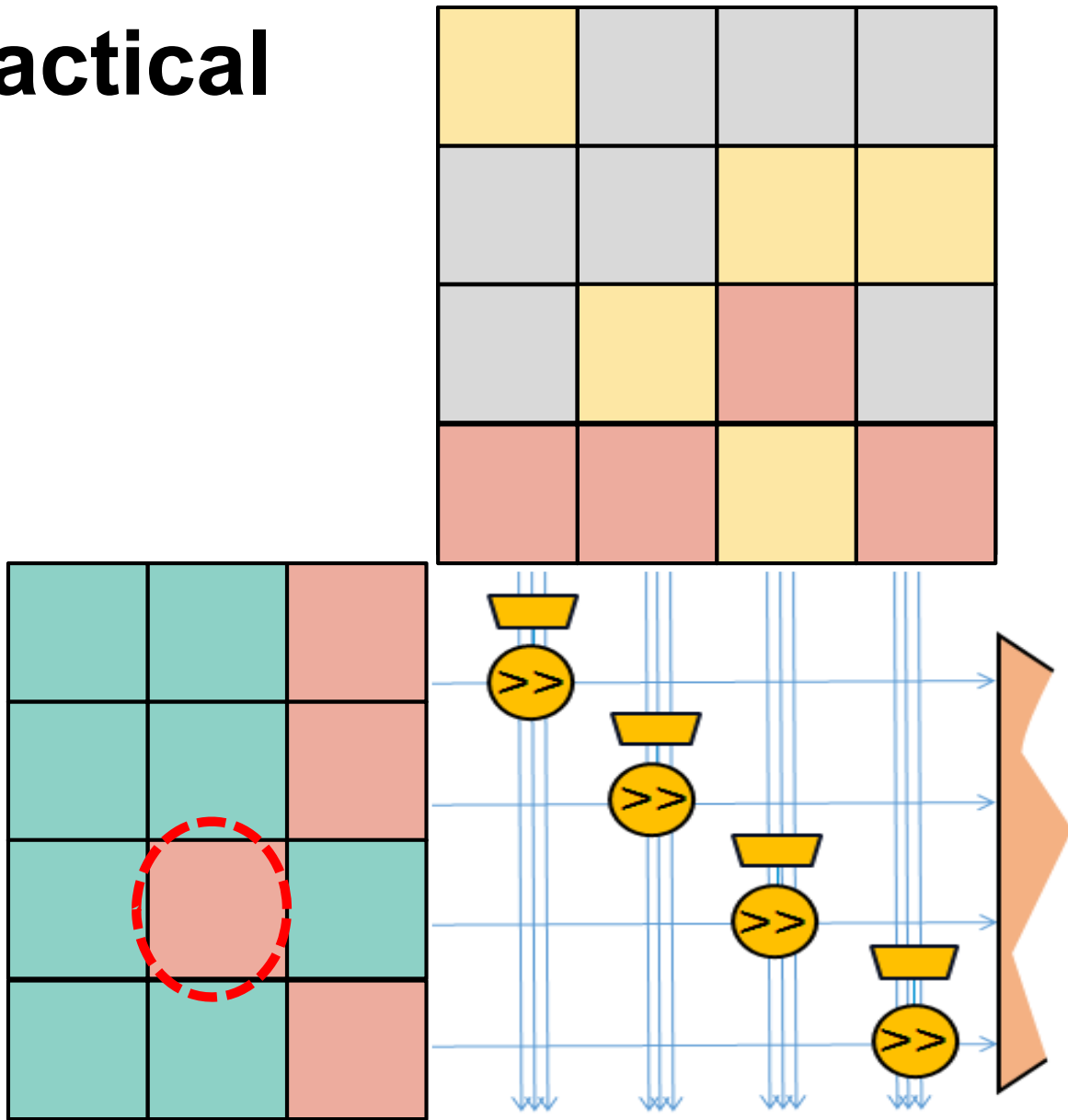
Do as you are told?

$$Out += A_0 \times W_0$$

$$Out += A_2 \times W_2$$

$$Out += 0 \times W_3$$

$$Out += A_4 \times W_4$$

.
.
.

Do as you are told?

$$\text{Out} += A_0 \times W_0$$

$$\text{Out} += A_2 \times W_2$$

$$\text{Out} += A_3 \times W_3$$

$$\text{Out} += A_4 \times W_4$$

⋮

# Do as you are told?

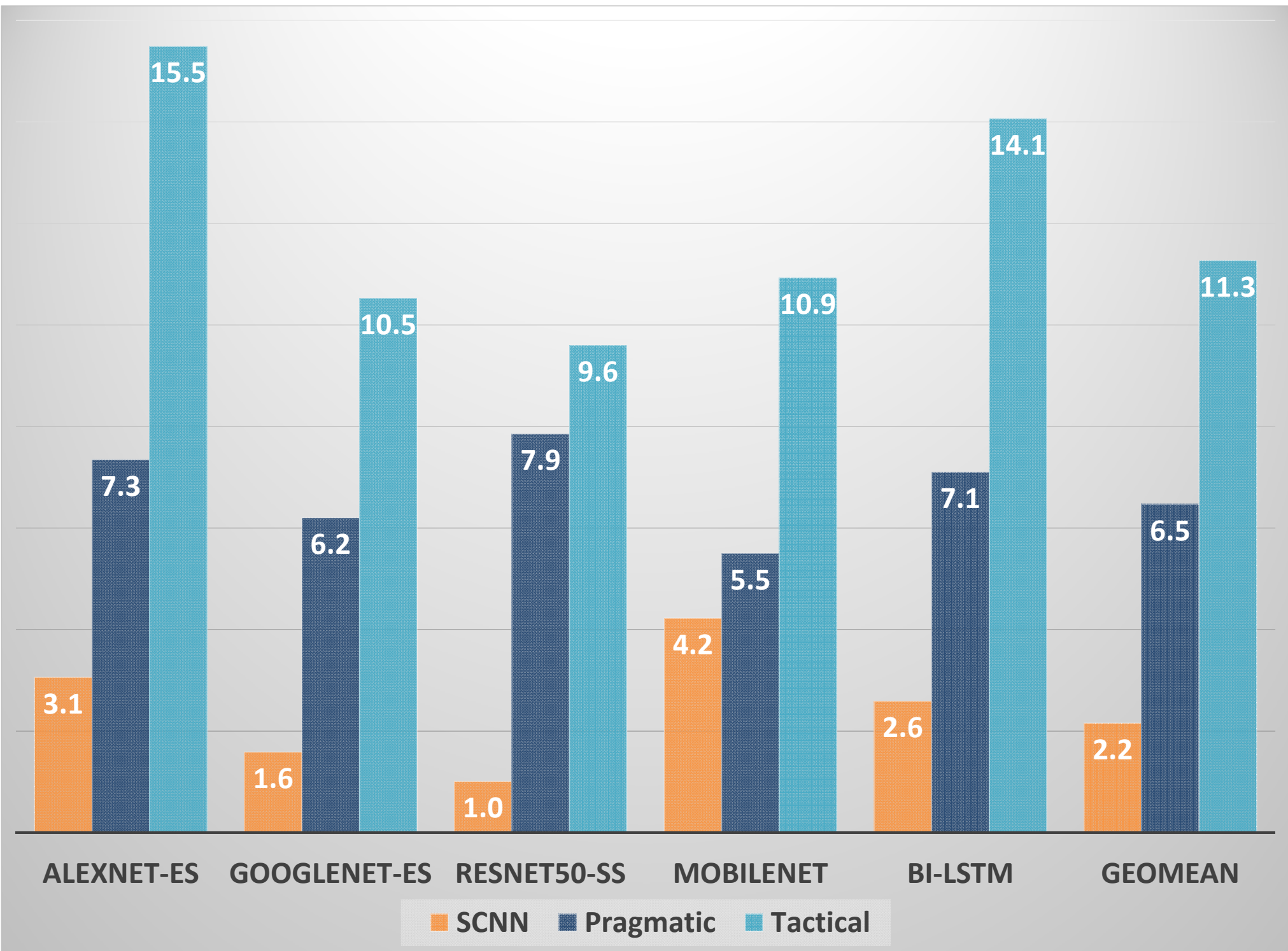Out += $000100100$ x $W_0$

Out += $000110100$ x $W_2$

Out += $000010000$ x $W_3$

Out += $010000010$ x $W_4$

⁝

# Do as you are told?

Out += $000100100$ x $W_0$

Out += $000110100$ x $W_2$

Out += $000010000$ x $W_3$

Out += $010000010$ x $W_4$

.
.
.

**Do as you are told?**

Out += $\quad\quad$ 1 $\quad$ 1 $\quad$ x $W_0$

Out += $\quad\quad$ 11 $\quad$ 1 $\quad$ x $W_2$

Out += $\quad\quad\quad$ 1 $\quad\quad\quad$ x $W_3$

Out += $\quad$ 1 $\quad\quad\quad$ 1 $\quad$ x $W_4$

# Making Typical Values Matter



*DPRed: Making Typical Activation and Weight Values Matter In Deep Learning Computing*, Delmas et al**.,**
**https://arxiv.org/abs/1804.06732**

# Fine-Grain Precision Adaptation

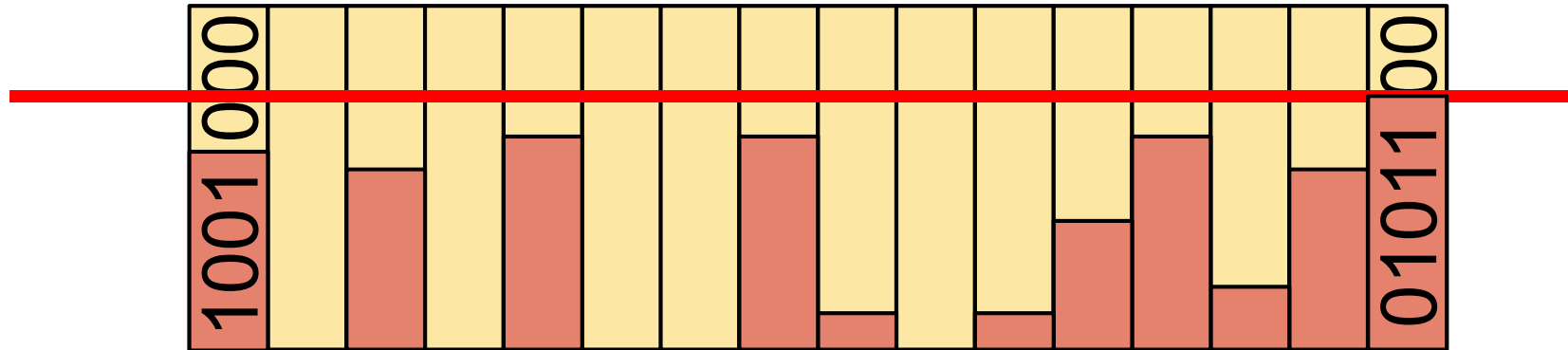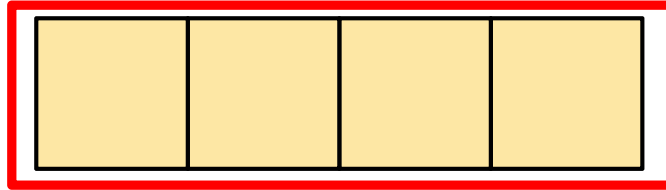**Detect Precisions:** On-the-fly for Activations + Statically for Weights



**In Memory**

| Zero Value Vector | | Prec | $V_0$ | $V_1$ | | $V_N$ | u |
|---|---|---|---|---|---|---|---|

16 — 4 — P — P — P

37

# Takeaways

Don't go after all W = 0

    Go for enough

    Software Scheduler + Restricted Motions are enough

Don't go after A = 0

    Go for bit sparsity 50% vs 90%+

Don't let the loud values dominate the data type

    Encode in groups

Bit-Tactical: Exploiting Ineffectual Computations in Convolutional Neural Networks: Which, Why, and How, Alberto Delmas, Patrick Judd, Dylan Malone Stuart, Zissis Poulos, Mostafa Mahmoud, Sayeh Sharify, Milos Nikolic, Andreas Moshovos, arXiv:1803.03688

```
W 0001 0100
A 0010 1010
```
_____

```
  0000 0000

  0010 1010

 0000 0000

 0010 1010

0000 0000

0010 1010

0000 0000

0000 0000
```
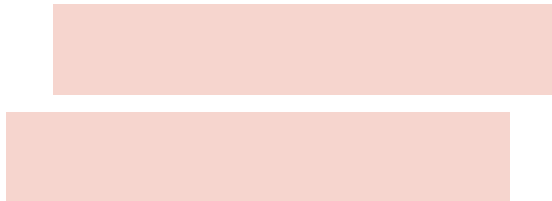
```
W 0001 0100
A 0010 1010
───────────────────

    0010 1010
    0000 0000
    0010 1010
    0000 0000
    0010 1010
```

```
W 0001 0100
A 0010 1010
```

0010 1010

0010 1010

0010 1010

```
W  0001  0100
A  0010  1010
```

```
   1   1  1
```

```
  1    1  1
```
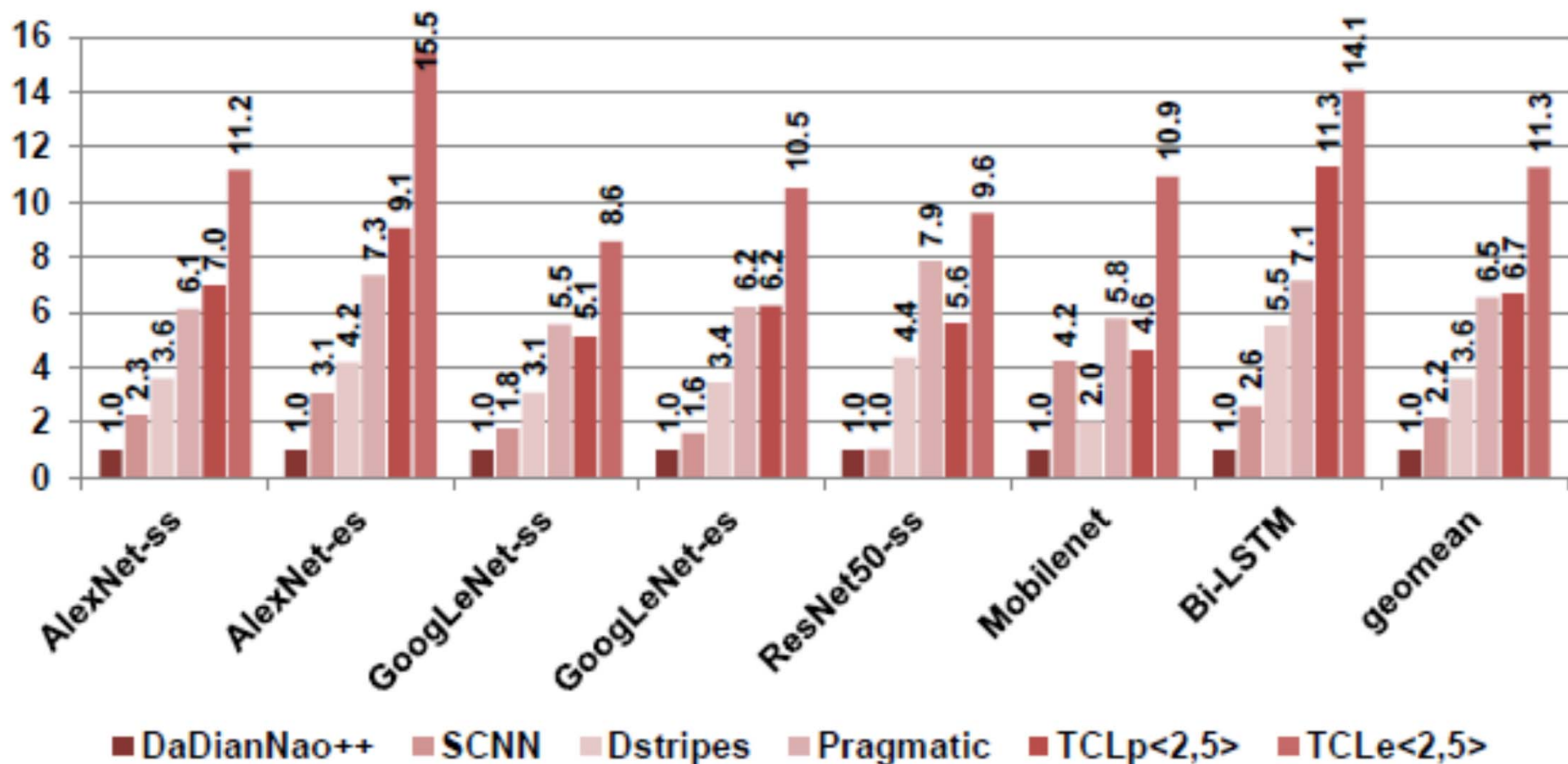
```
 1    1  1
```

# Laconic

- **Best for edge devices**

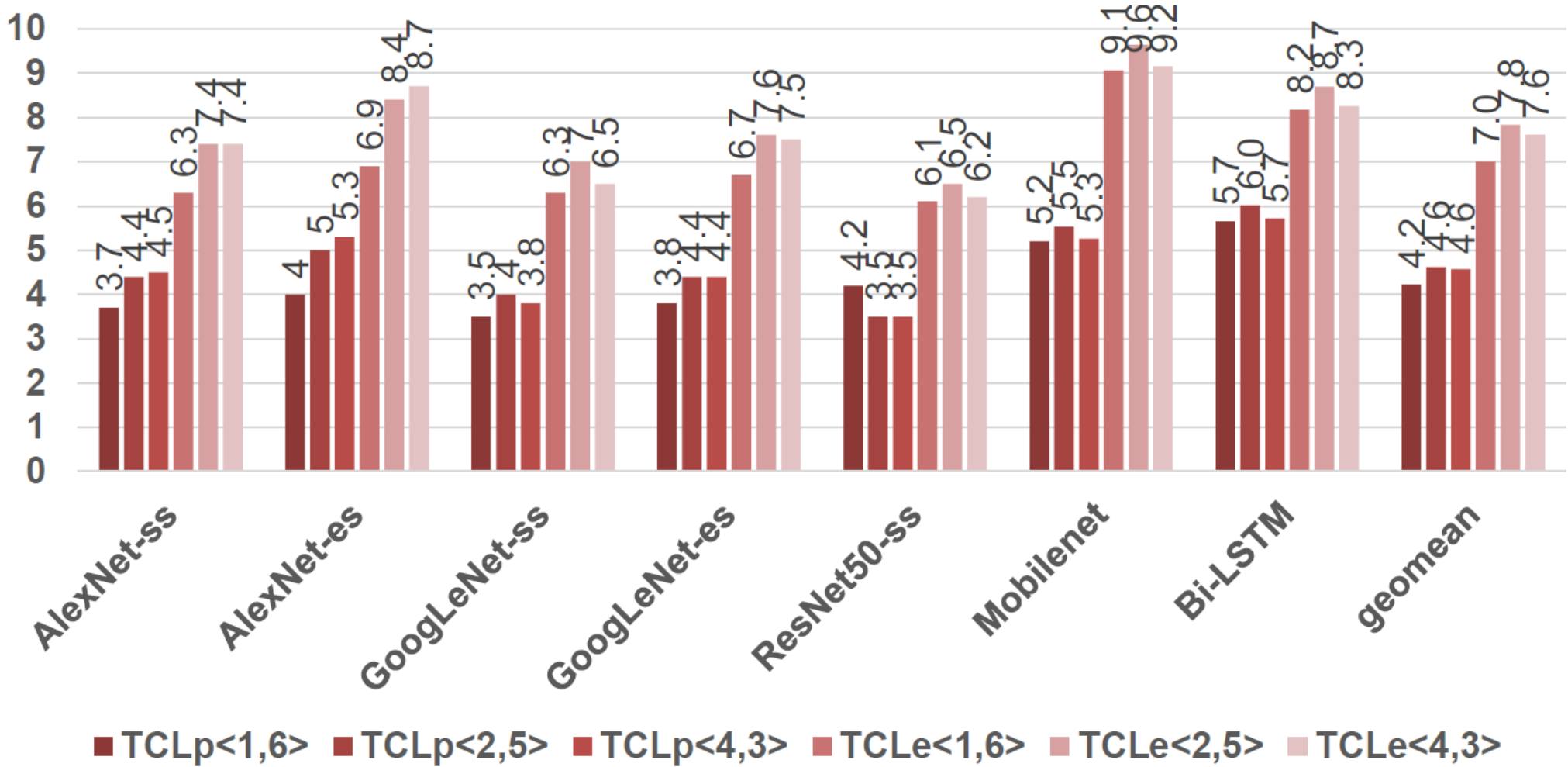- **Laconic Deep Learning Computing,** *Sayeh Sharify, Mostafa Mahmoud, Alberto Delmas Lascorz, Milos Nikolic, Andreas Moshovos, Arxiv,* **arXiv:1805.04513**

Figure 12: Speedup with 8b quantization