# pCNN: Parallel Convolutional Neural Network Implementations for Handwritten Digit Recognition
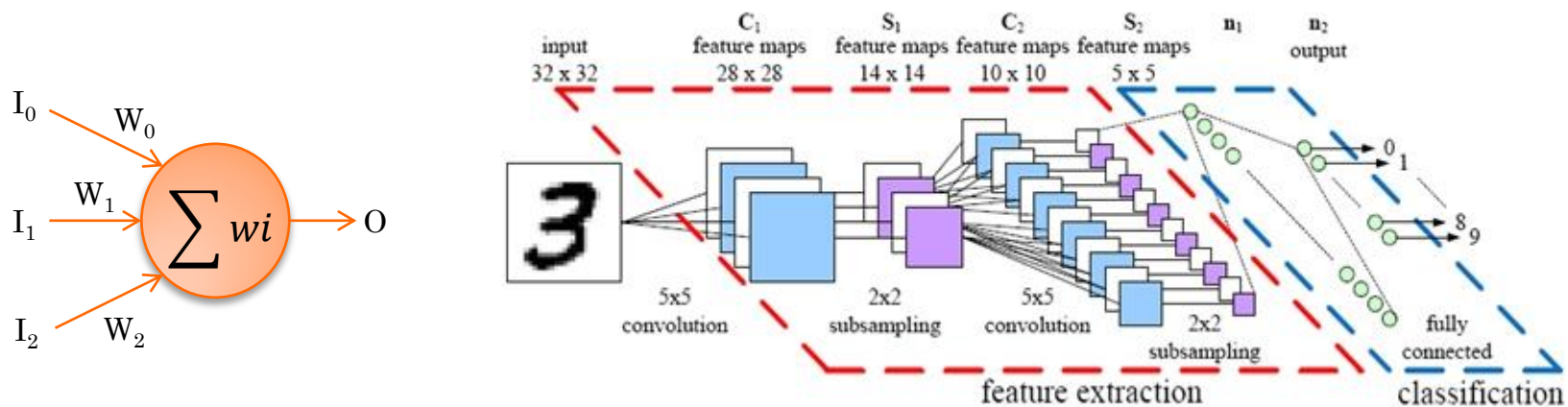
**Swagath Venkataramani**

**Rangharajan Venkatesan**

**Ashiwan Sivakumar**

1

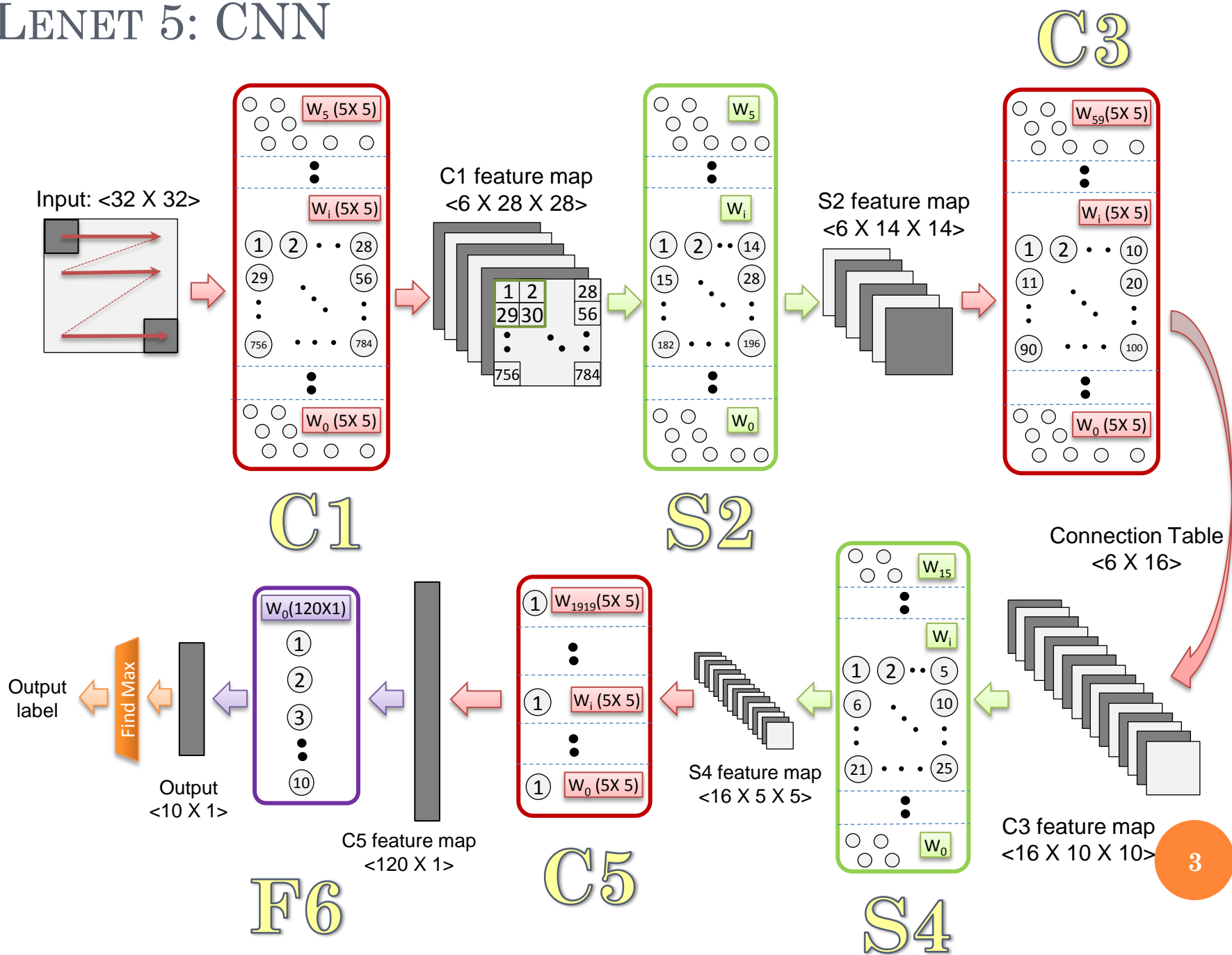# OBJECTIVE

- Implement parallel software versions of Convolutional Neural Networks (CNN)
  - OpenMP, pthreads, MPI
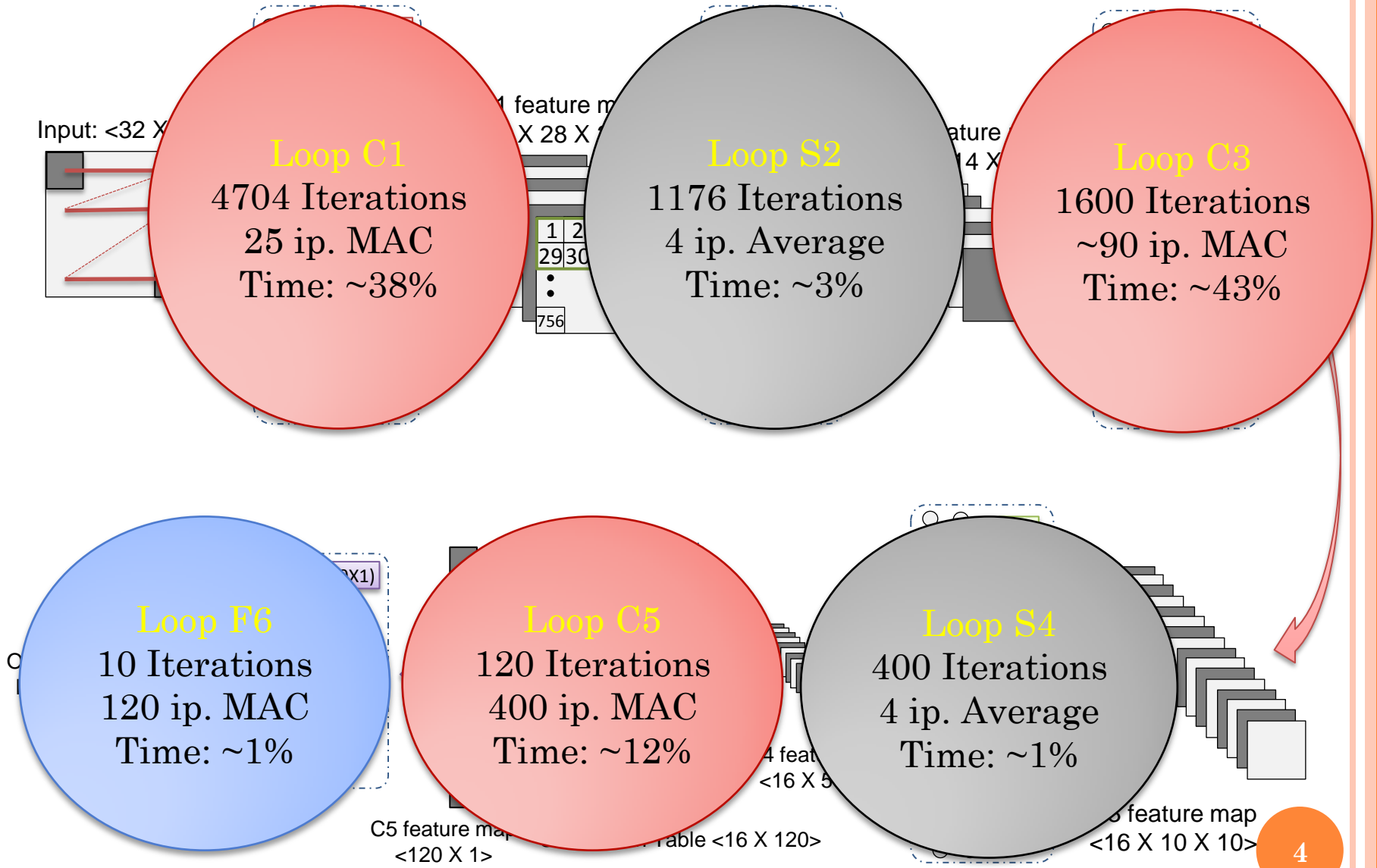- Lenet-5: Designed for handwritten digit recognition application.



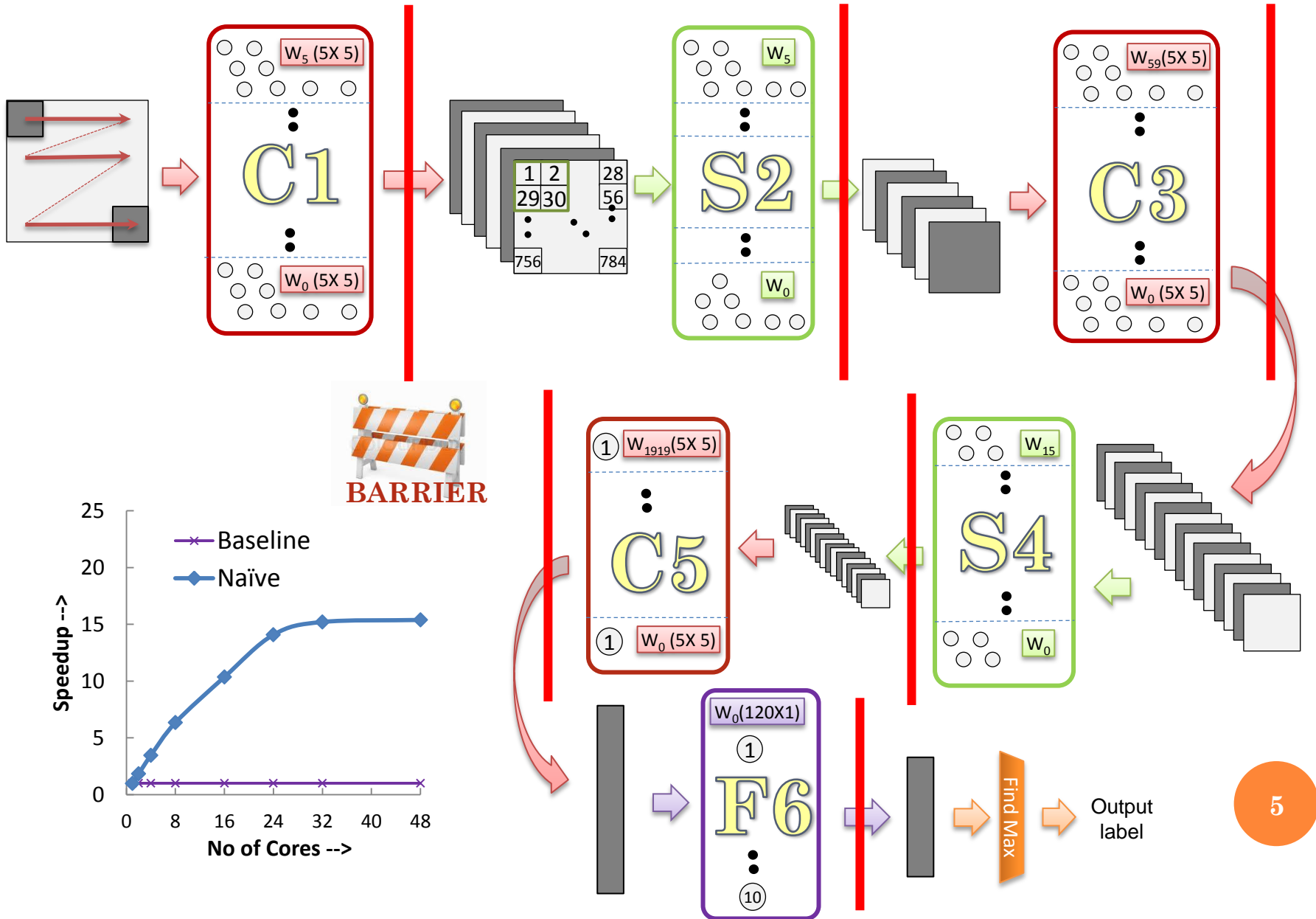- MNIST dataset
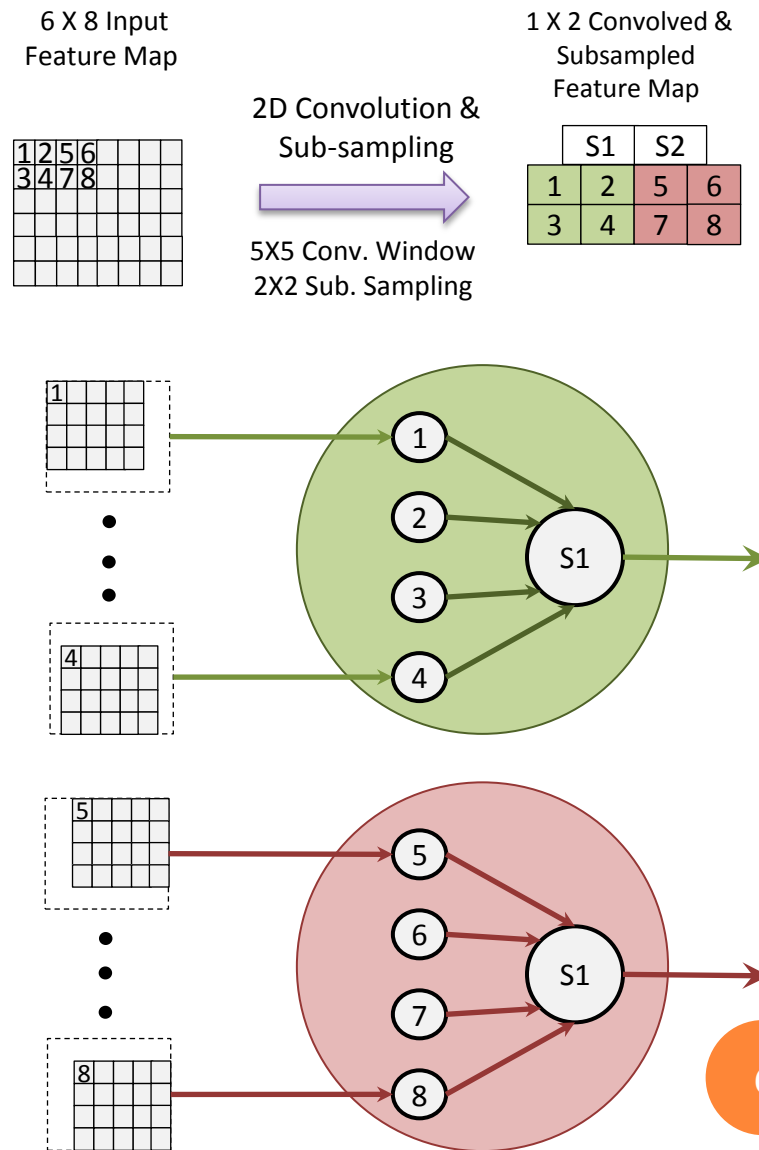- Experiments run on server with 48 AMD cores
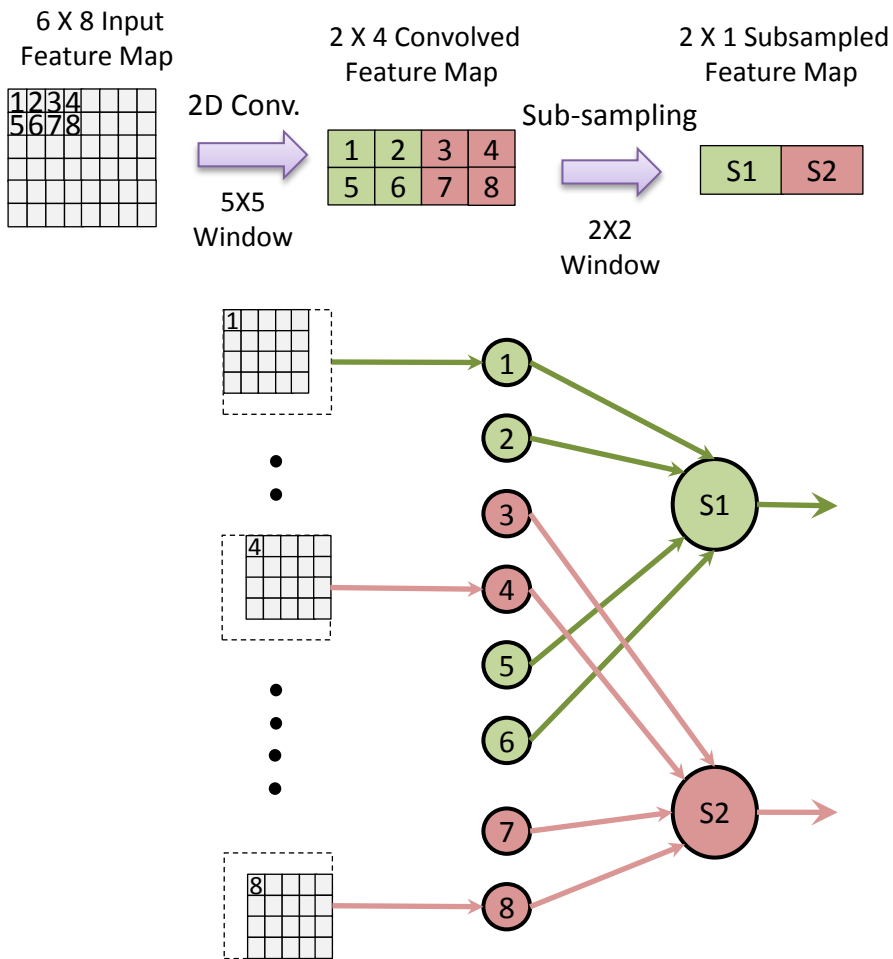
2

# LENET 5: CNN

Input: <32 X 32>

**C1**

$W_5$ (5X 5)

$W_i$ (5X 5)

① ② • • ㉘
㉙         ㊱
⑦⑤⑥    ⑦⑧④

$W_0$ (5X 5)

C1 feature map
<6 X 28 X 28>

| 1 | 2 | | 28 |
|---|---|---|---|
| 29 | 30 | | 56 |

756      784

**S2**

$W_5$

$W_i$

① ② • • ⑭
⑮         ㉘
⑧②  • • • ⑨⑥

$W_0$

S2 feature map
<6 X 14 X 14>

**C3**

$W_{59}$(5X 5)

$W_i$ (5X 5)

① ② • • ⑩
⑪         ⑳
⑨⓪ • • • ⑩⓪

$W_0$ (5X 5)

Connection Table
<6 X 16>

C3 feature map
<16 X 10 X 10>

**S4**

$W_{15}$

$W_i$

① ② • • ⑤
⑥         ⑩
㉑ • • • ㉕

$W_0$

S4 feature map
<16 X 5 X 5>

**C5**

① $W_{1919}$(5X 5)

① $W_i$ (5X 5)

① $W_0$ (5X 5)

C5 feature map
<120 X 1>

**F6**

$W_0$(120X1)

①
②
③
•
⑩

Output
<10 X 1>

Find Max

Output
label

**3**

Input: <32 X

1 feature m
X 28 X

ature
14 X

| 1 | 2 |
|---|---|
| 29 | 30 |

756

**Loop C1**
4704 Iterations
25 ip. MAC
Time: ~38%

**Loop S2**
1176 Iterations
4 ip. Average
Time: ~3%

**Loop C3**
1600 Iterations
~90 ip. MAC
Time: ~43%

**Loop F6**
10 Iterations
120 ip. MAC
Time: ~1%

**Loop C5**
120 Iterations
400 ip. MAC
Time: ~12%

**Loop S4**
400 Iterations
4 ip. Average
Time: ~1%

C5 feature map
<120 X 1>

Table <16 X 120>

4 feat
<16 X 5

feature map
<16 X 10 X 10>

4

6 X 8 Input
Feature Map

2 X 4 Convolved
Feature Map

2 X 1 Subsampled
Feature Map

2D Conv.

| 1 | 2 | 3 | 4 |
| 5 | 6 | 7 | 8 |

Sub-sampling

5X5
Window

| S1 | S2 |

2X2
Window

6 X 8 Input
Feature Map

2D Convolution &
Sub-sampling

1 X 2 Convolved &
Subsampled
Feature Map

5X5 Conv. Window
2X2 Sub. Sampling

| | S1 | | S2 | |
| 1 | 2 | 5 | 6 |
| 3 | 4 | 7 | 8 |

6

$W_5$ (5X 5) $W_5$

$W_0$ (5X 5) $W_0$

C1S2

$W_{59}$(5X 5) $W_{59}$

$W_0$ (5X 5) $W_0$

C3S4

① $W_{1919}$(5X 5)

① $W_0$ (5X 5)

C5

$W_0$(120X1)

①

F6
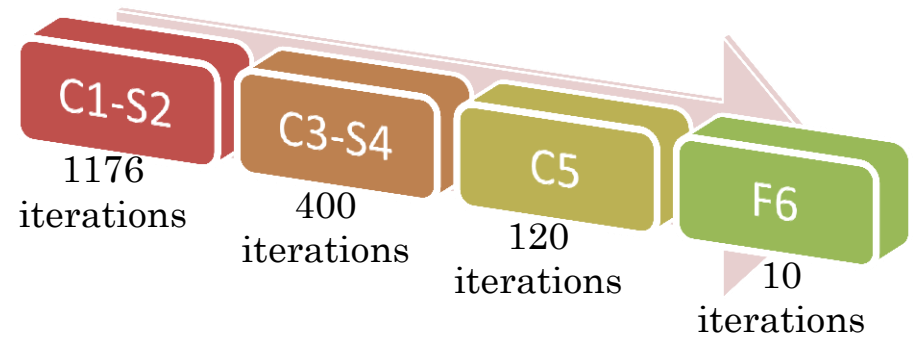
⑩

Find Max

Output label

Speedup -->

No of Cores -->

Baseline
Naïve
Fused

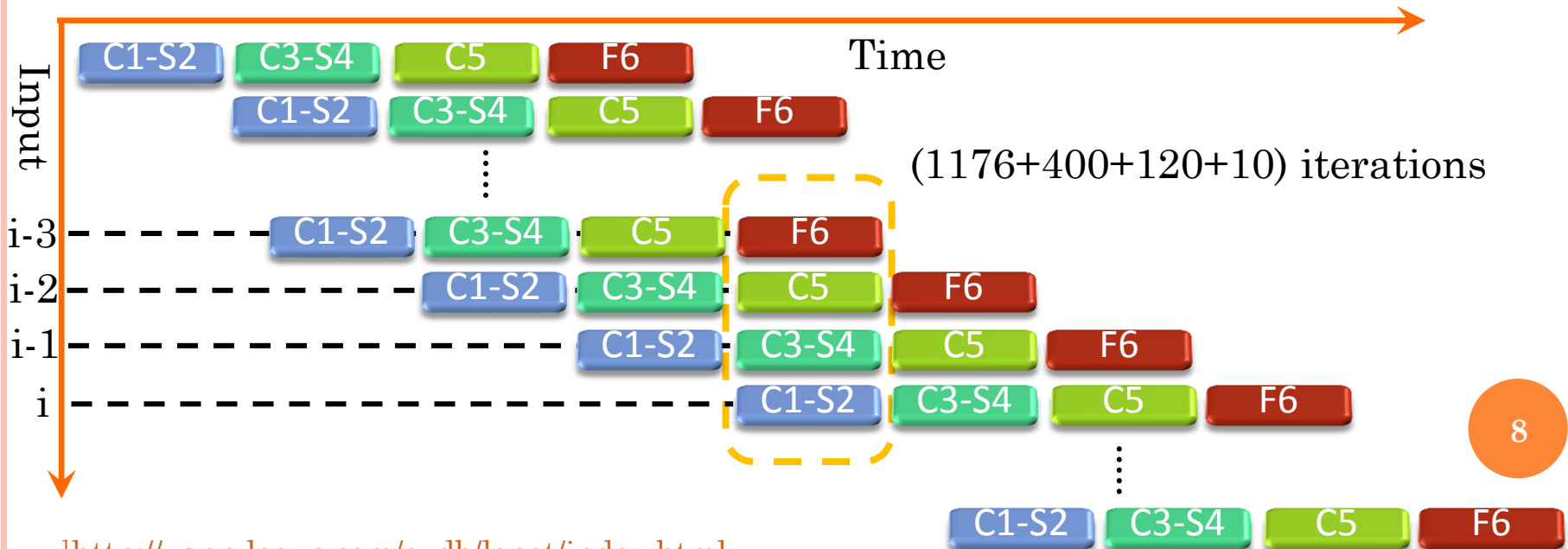# TRANSFORMATION 3: PIPE-FUSED PARALLEL



Digit recognition typically processes stream of i/ps[1]



Producer-Consumer relationship across layers in Fused Parallel implementation
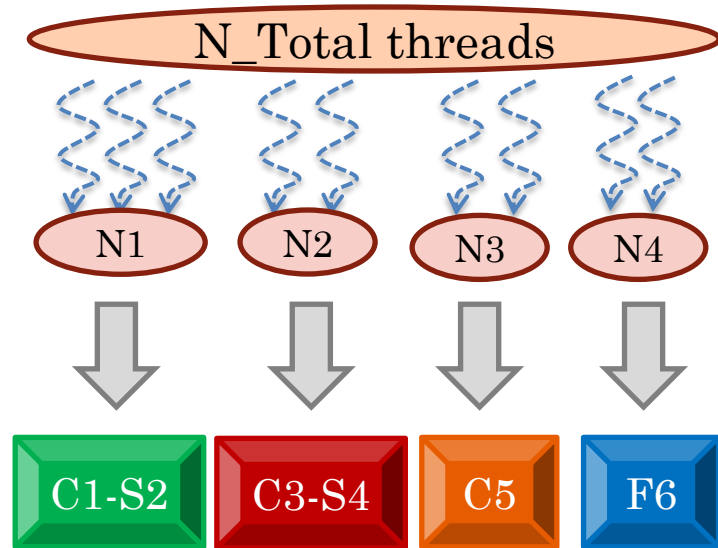
## PIPE- FUSED: Enhanced parallelism through Pipelining



(1176+400+120+10) iterations

8

Pseudo-code

```
#pragma omp for
for (i=1:N_Total)
if ( i < N1)
        process C1-S2
else if ( i< N1+N2)
        process C3-S4
else if ( i < N1+N2+N3)
        process C5
else
        process F6
```
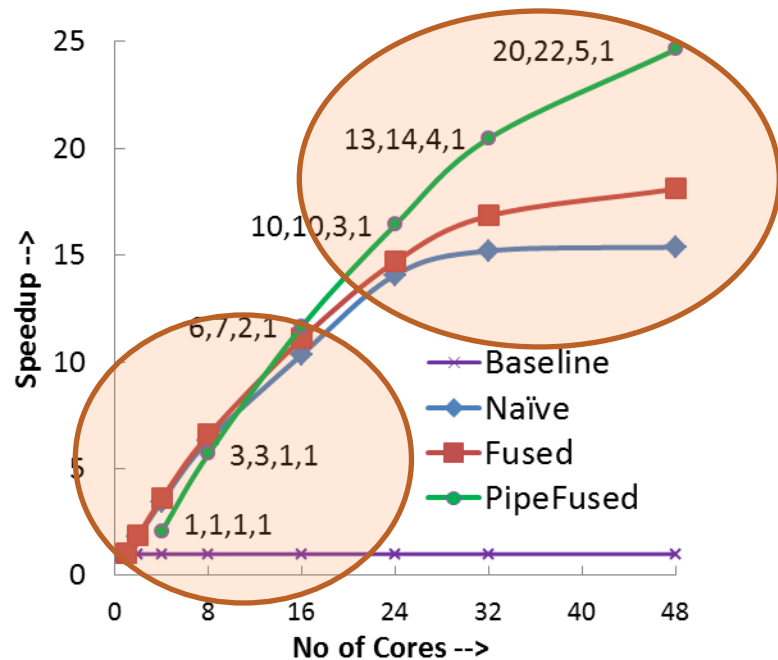


- **N1,N2,N3,N4??**
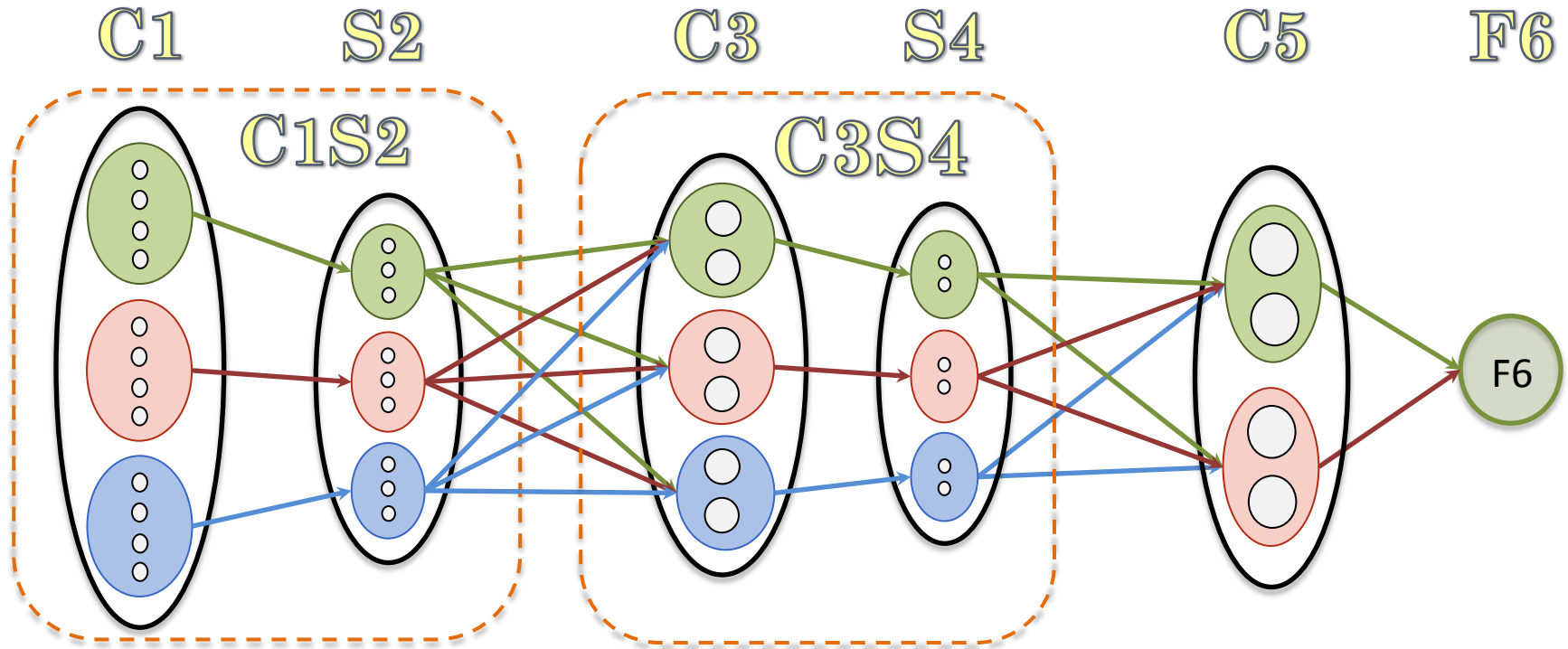
- **Design-space exploration**

# TRANSFORMATION 3: PIPEFUSED PARALLEL

Pseudo-code

```
#pragma omp for
for (i=1:N_Total)
if ( i < N1)
        process C1-S2
else if ( i< N1+N2)
        process C3-S4
else if ( I < N1+N2+N3)
        process C5
else
        process F6
```

# DISTRIBUTED MEMORY MODEL - MPI



Naive implementation:
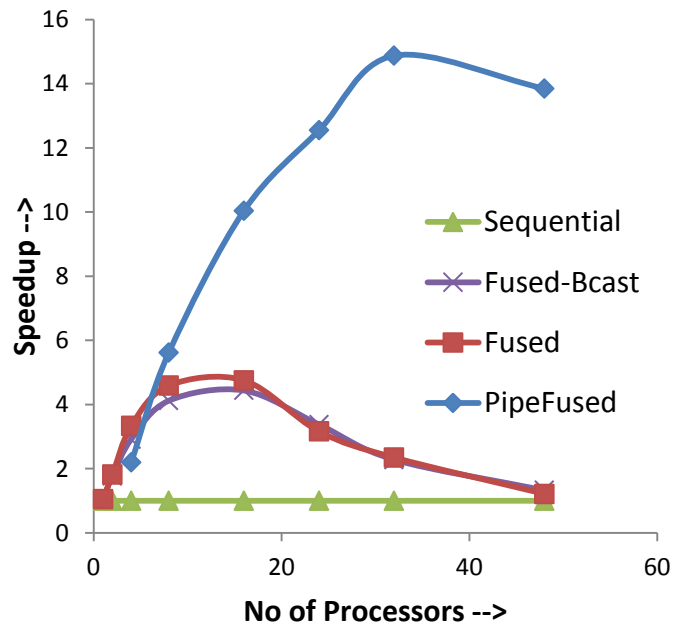
➢ Broadcast output before running next layer

Transformation 1: Fuse layers
➢ Eliminates C1-S2 and C3-S4 communication
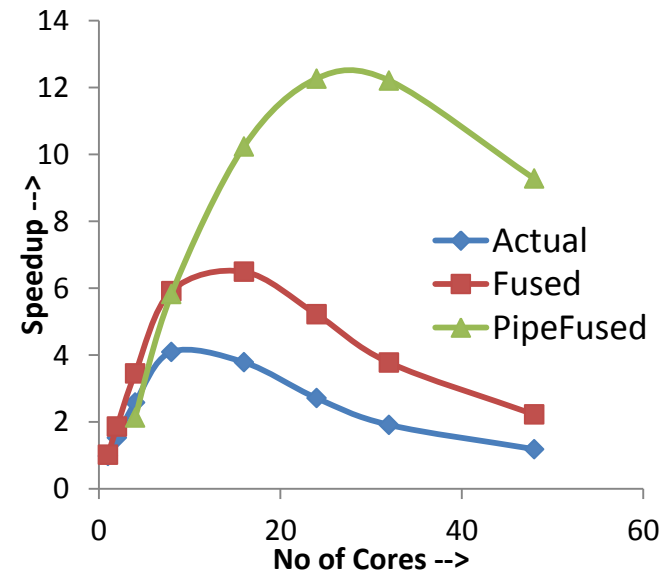➢ Still broadcast between S2-C3 and S4-C5

Transformation 2: "Selective Send" based on connection table

# RESULTS

# SUMMARY & FUTURE WORK

- Intense Communication between Neurons – Distributed memory model suffers

- Loop body of each neuron is small – Fork-Join overheads

- Take advantage of "Convolution followed by Sub-sampling"

- Pipe-fused expands the parallelism beyond each network layer

- Parallelize training phase
  - OpenMP and MPI

13