

FW  NXT

Efficient compiler code generation for  
Deep Learning Snowflake co-processor

# Deep Learning



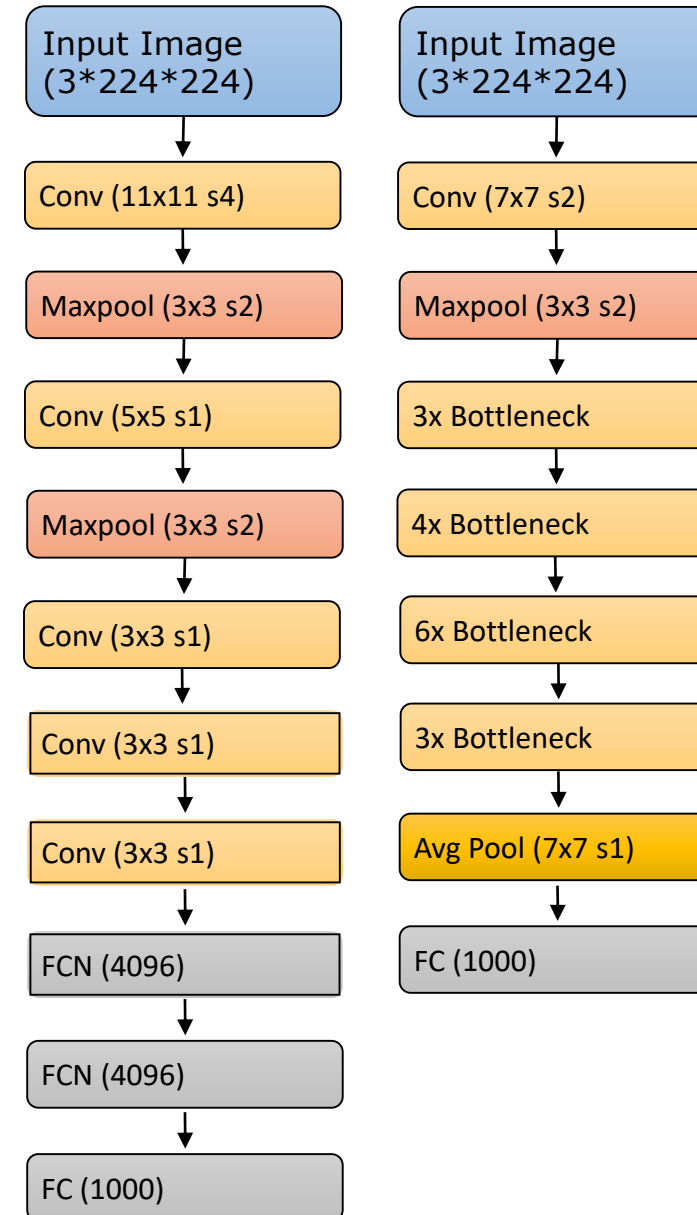
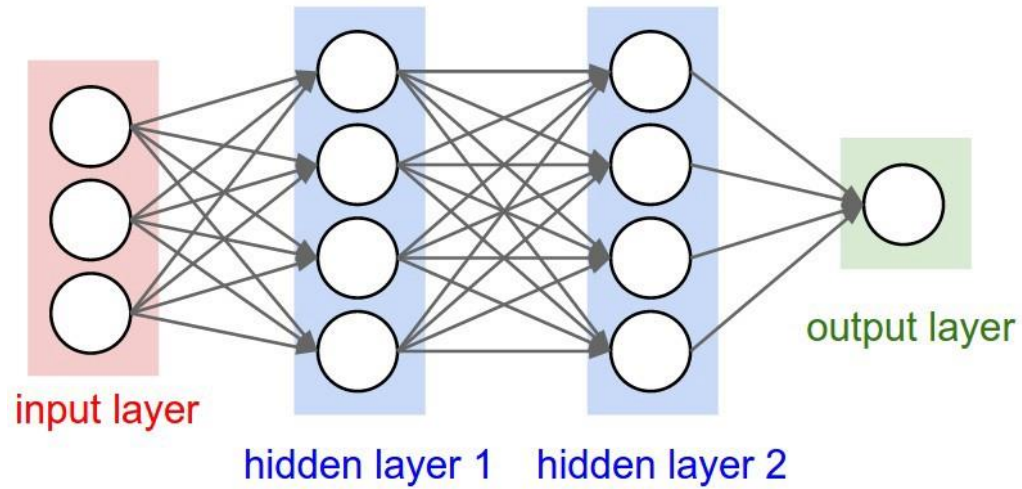


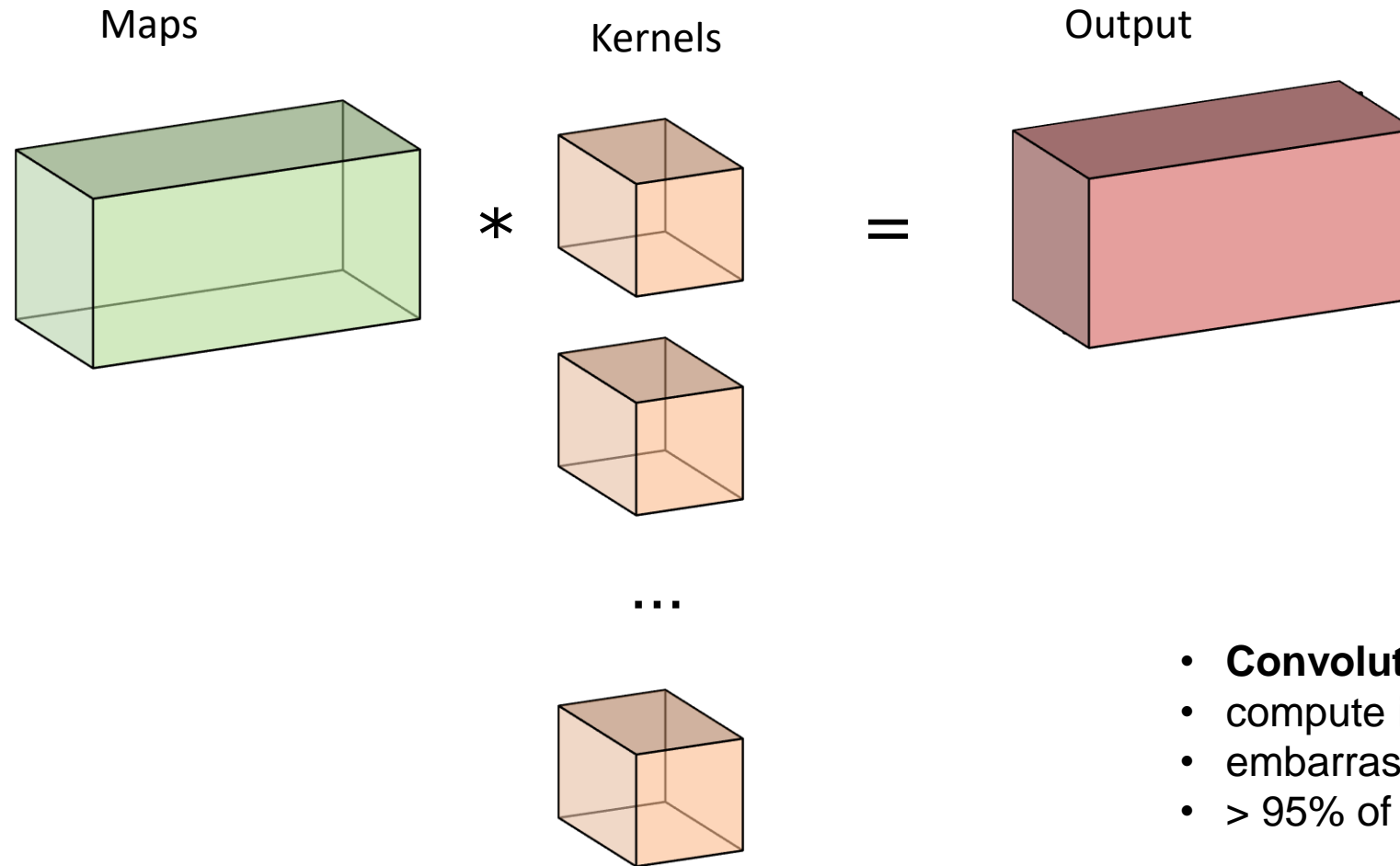
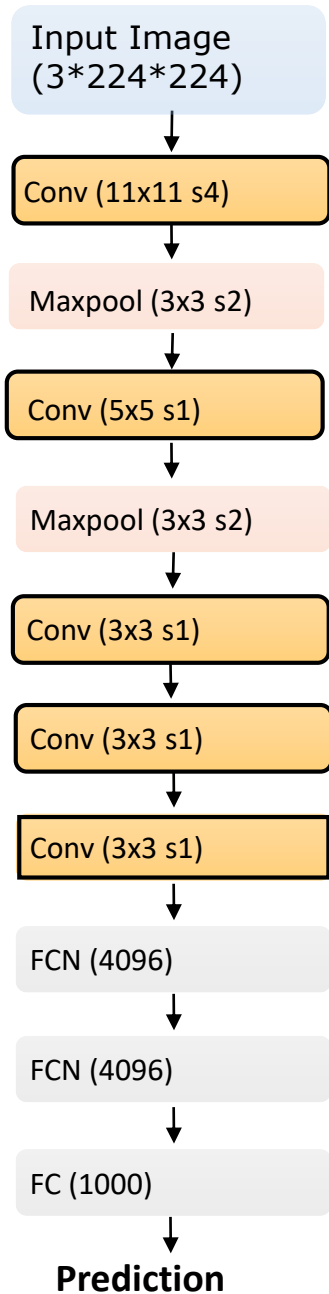






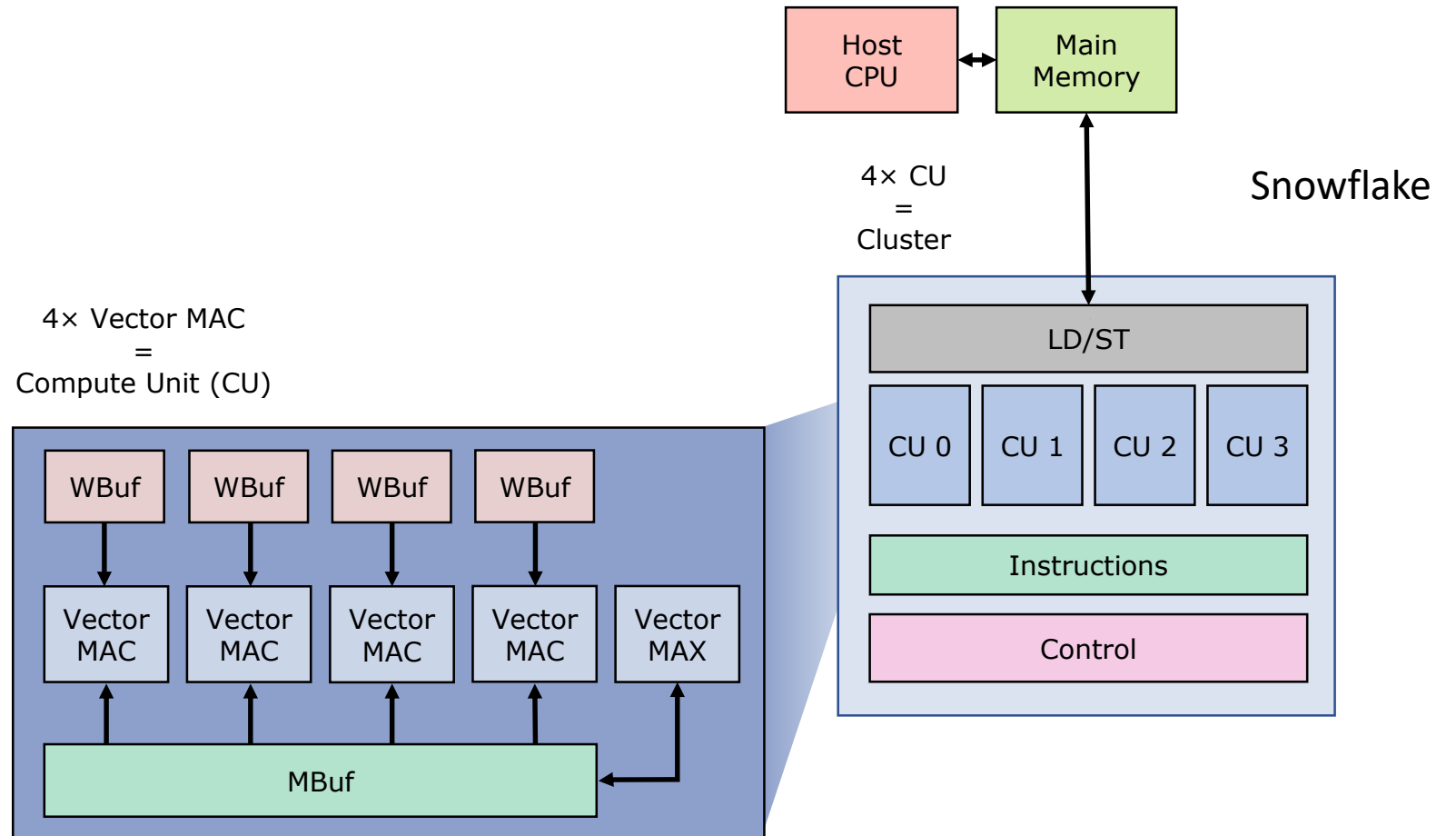
# Conv Neural Networks





- **Convolutions:**
- compute intensive
- embarrassingly parallel
- > 95% of the workload

# Co-processor





# Snowflake instruction set

Move

Add

Multiply

Branch

MAC : Multiply and accumulate a trace

MAX: Compare a trace

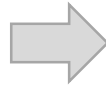
VMOV: Pre-load values

LOAD: mem2buffer

TMOV: buffer2mem

# Goal

```
1 import torch
2 import torch.nn as nn
3 class ModelDef(nn.Module):
4     def __init__(self, num_classes=1000):
5         super(ModelDef, self).__init__()
6         self.features = nn.Sequential(
7             nn.Conv2d(3, 64, kernel_size=11, stride=4, padding=2),
8             nn.ReLU(inplace=True),
9             nn.MaxPool2d(kernel_size=3, stride=2),
10            nn.Conv2d(64, 192, kernel_size=5, padding=2),
11            nn.ReLU(inplace=True),
12            nn.MaxPool2d(kernel_size=3, stride=2),
13            nn.Conv2d(192, 384, kernel_size=3, padding=1),
14            nn.ReLU(inplace=True),
15            nn.Conv2d(384, 256, kernel_size=3, padding=1),
16            nn.ReLU(inplace=True),
17            nn.Conv2d(256, 256, kernel_size=3, padding=1),
18            nn.ReLU(inplace=True),
19            nn.MaxPool2d(kernel_size=3, stride=2),
20        )
21        self.classifier = nn.Sequential(
22            nn.Dropout(),
23            nn.Conv2d(256, 4096, kernel_size=6),
24            nn.ReLU(inplace=True),
25            nn.Dropout(),
26            nn.Conv2d(4096, 4096, kernel_size=1),
27            nn.ReLU(inplace=True),
28            nn.Conv2d(4096, num_classes, kernel_size=1),
29            # nn.Softmax(),
30        )
31    def forward(self, x):
32        x = self.features(x)
33        x = self.classifier(x)
34        return x
```

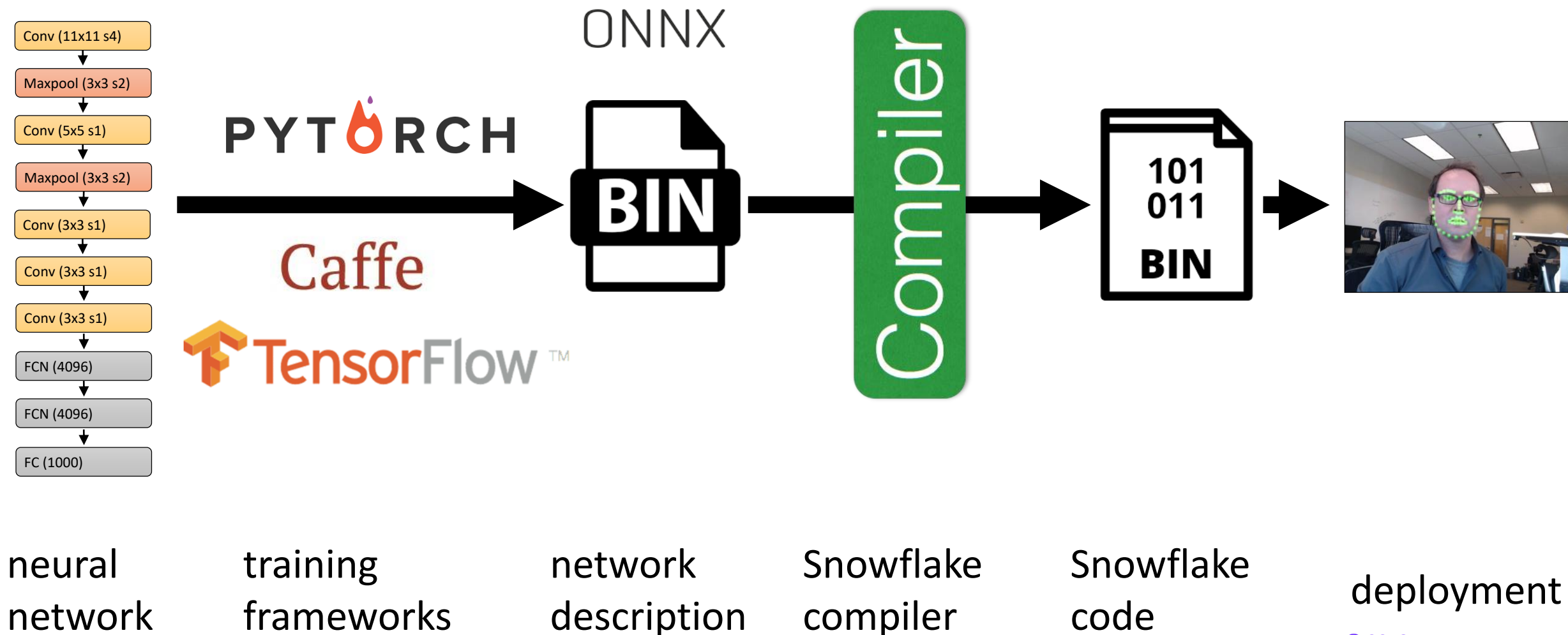


```
comp_add(0,1,22,0,0x000000f);
data_mov(1,18,1,0x0000000);
comp_add(0,0,9,0,0x0000000);
data_mov(0,3,0,0x3ffffff);
data_mov(0,2,0,0x3ffffff);
comp_mac(1,0x1f,0,1,0x000300);
brch_inst('l',21,2,0x000fff);
comp_add(0,0,0,0,0x001380);
comp_add(0,1,1,0,0x004800);
```

...

- Correctness
- Speed
- Extendibility

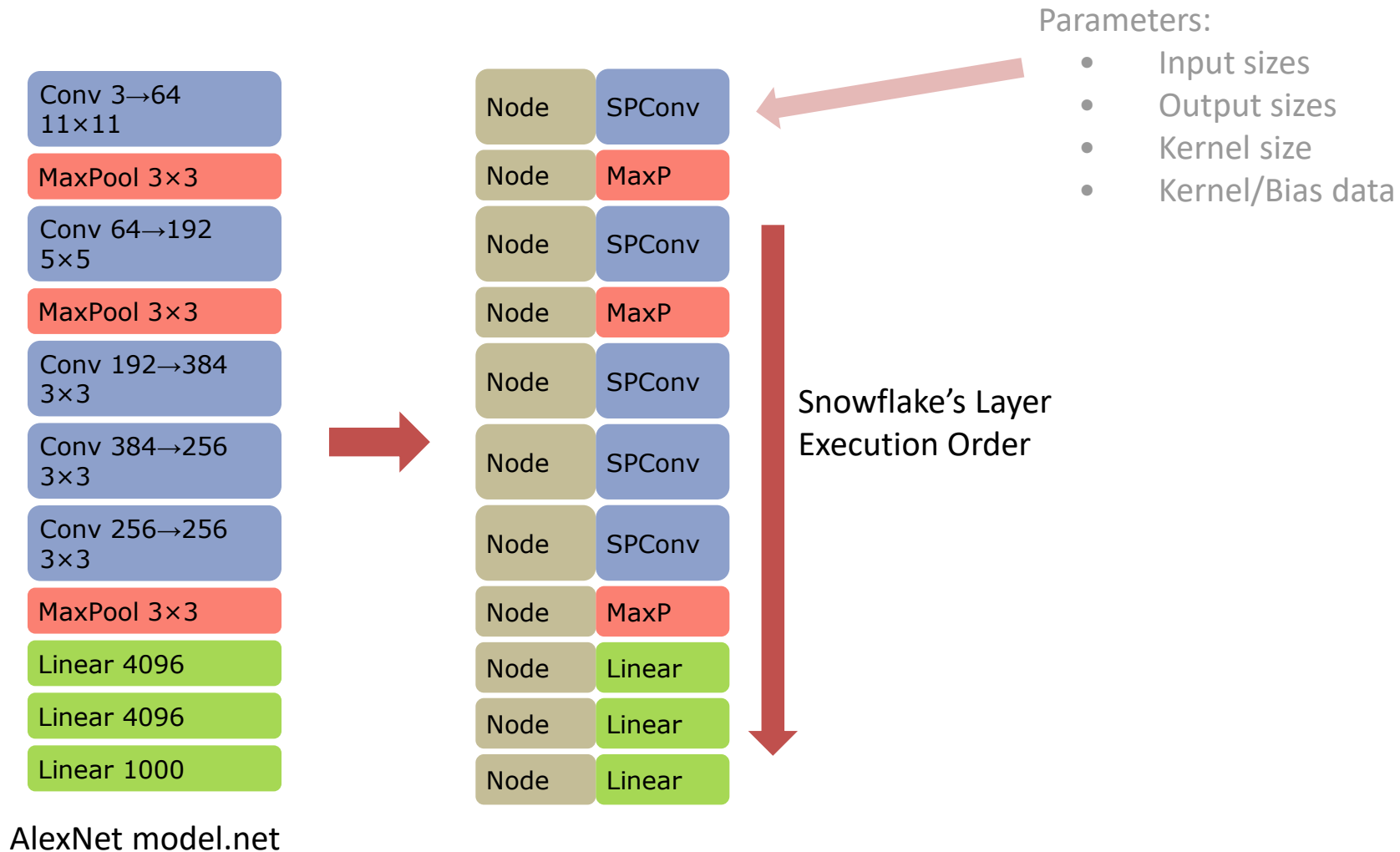
# Overview



# Snowflake compiler steps

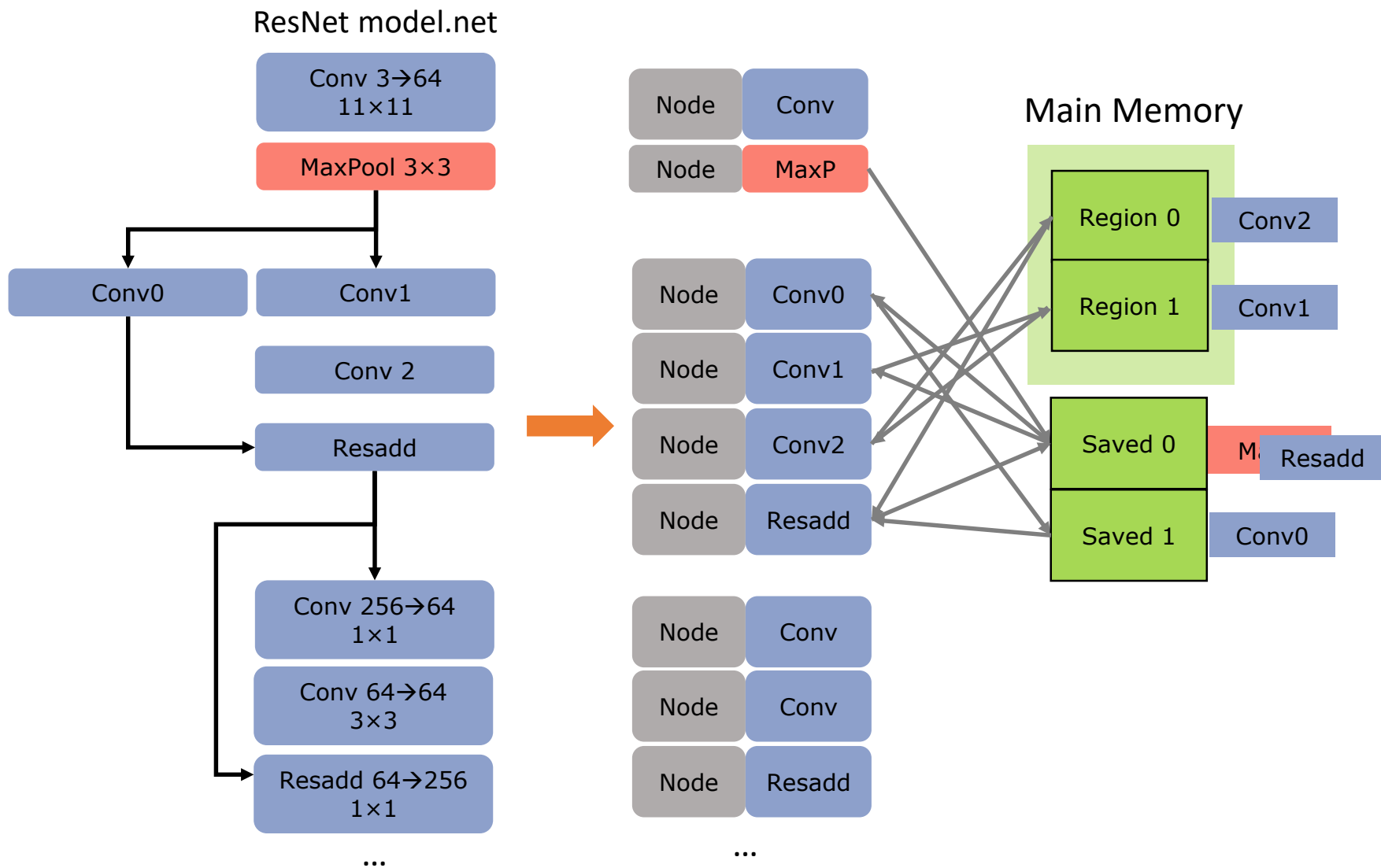
- Parsing
- Partition
- Code generation

# Parse

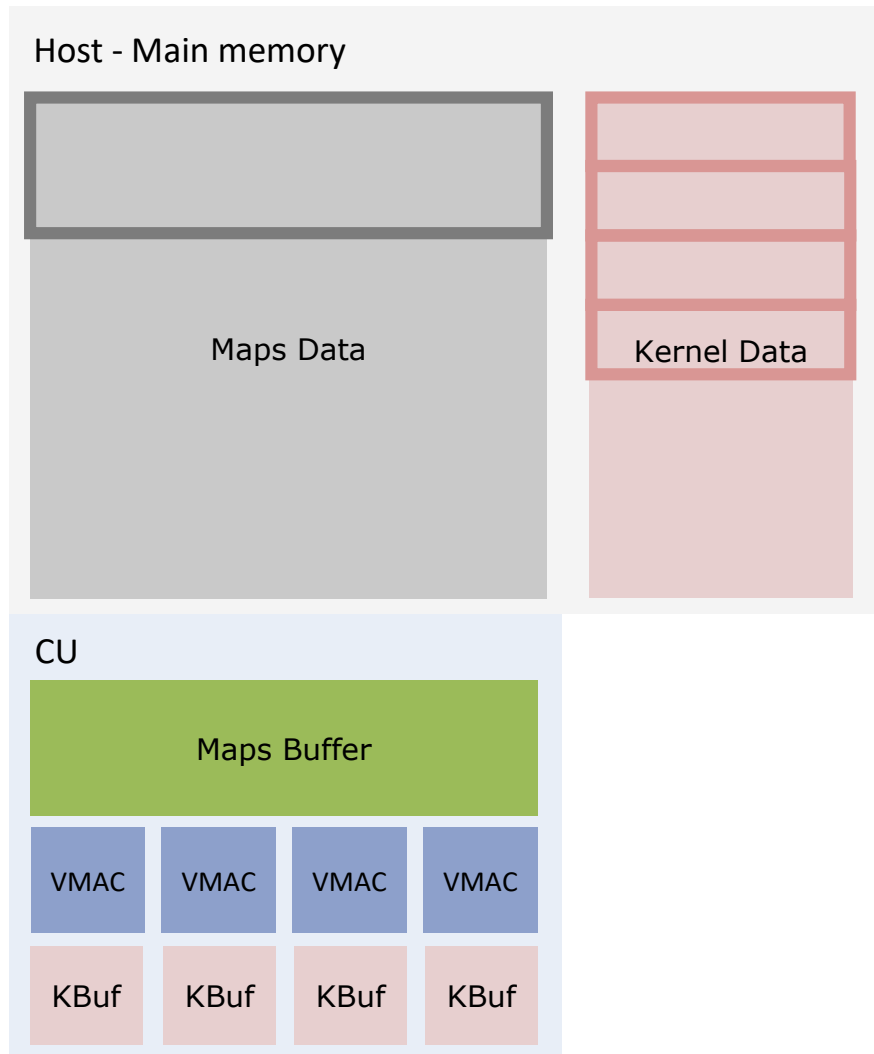


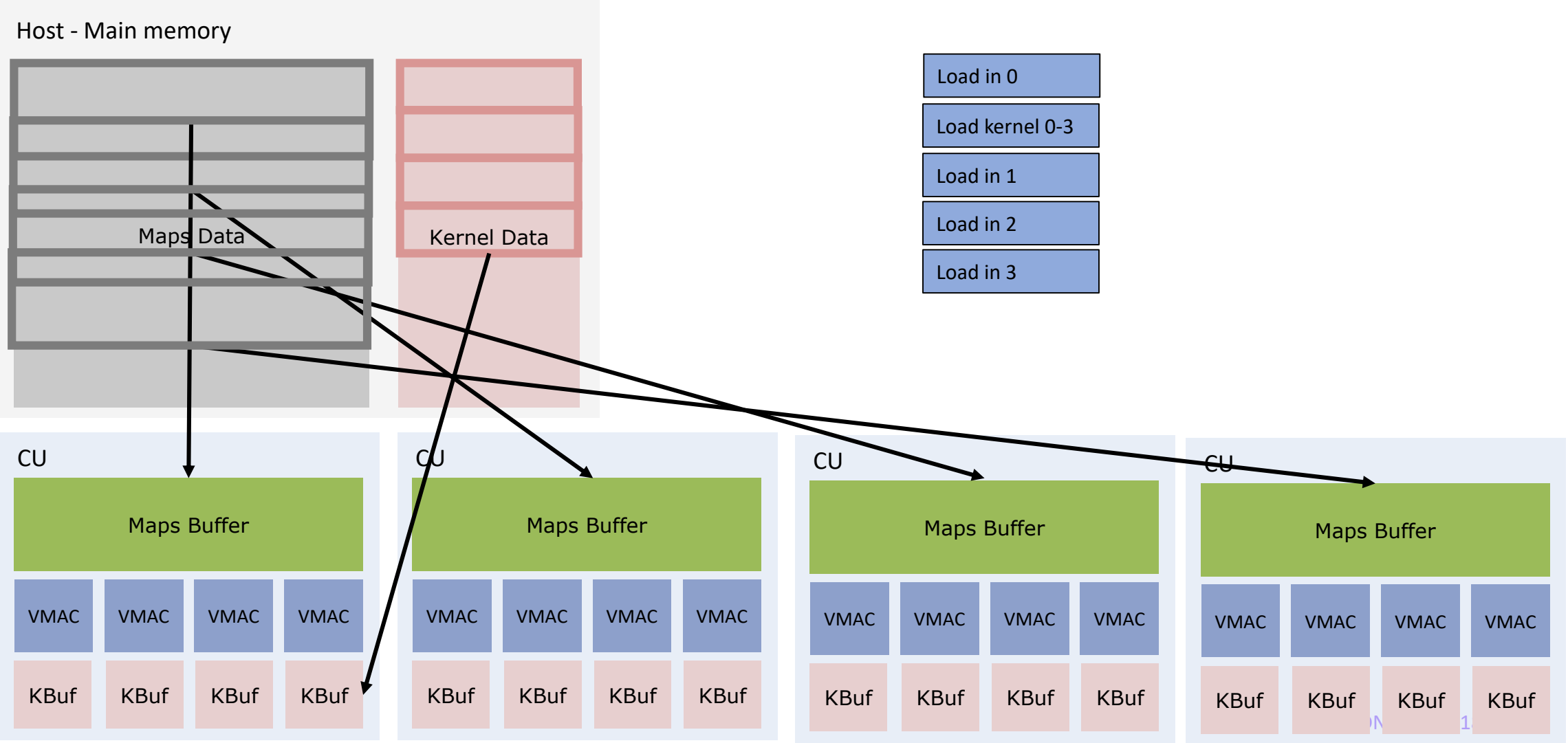






# Partition





# Order

maps

m(0)

m(1)

m(2)

kernel

k(0)

k(1)

k(2)

```
for( k=0;k<3;k++ )  
  for ( m=0;m<3;m++ )
```

k(0)

m(0)

k(0)

m(1)

k(0)

m(2)

k(1)

m(0)

k(1)

m(1)

k(1)

m(2)

...

```
for ( m=0;m<3;m++ )  
  for( k=0;k<3;k++ )
```

k(0)

m(0)

k(1)

m(0)

k(2)

m(0)

k(0)

m(1)

k(1)

m(1)

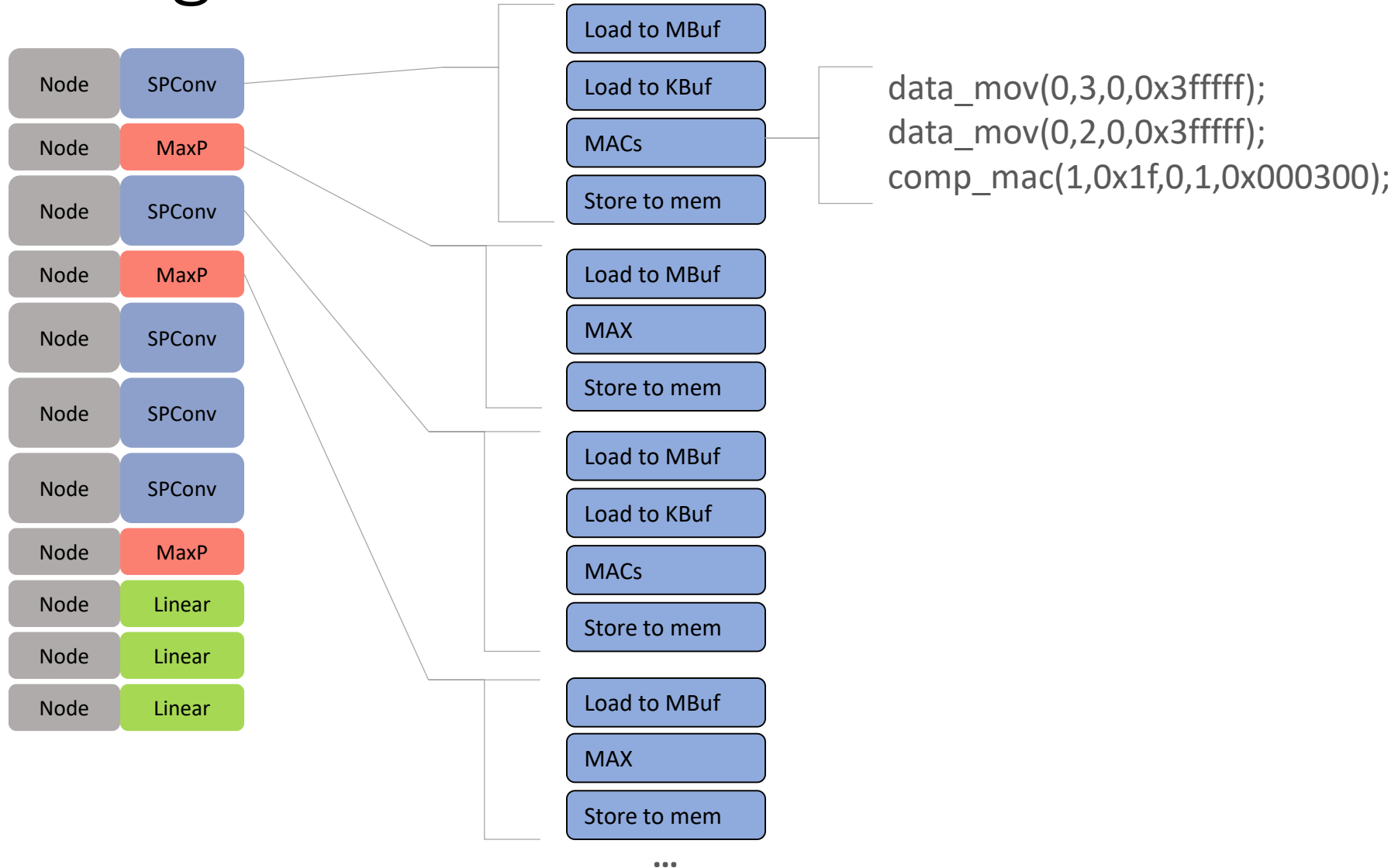
k(1)

m(1)

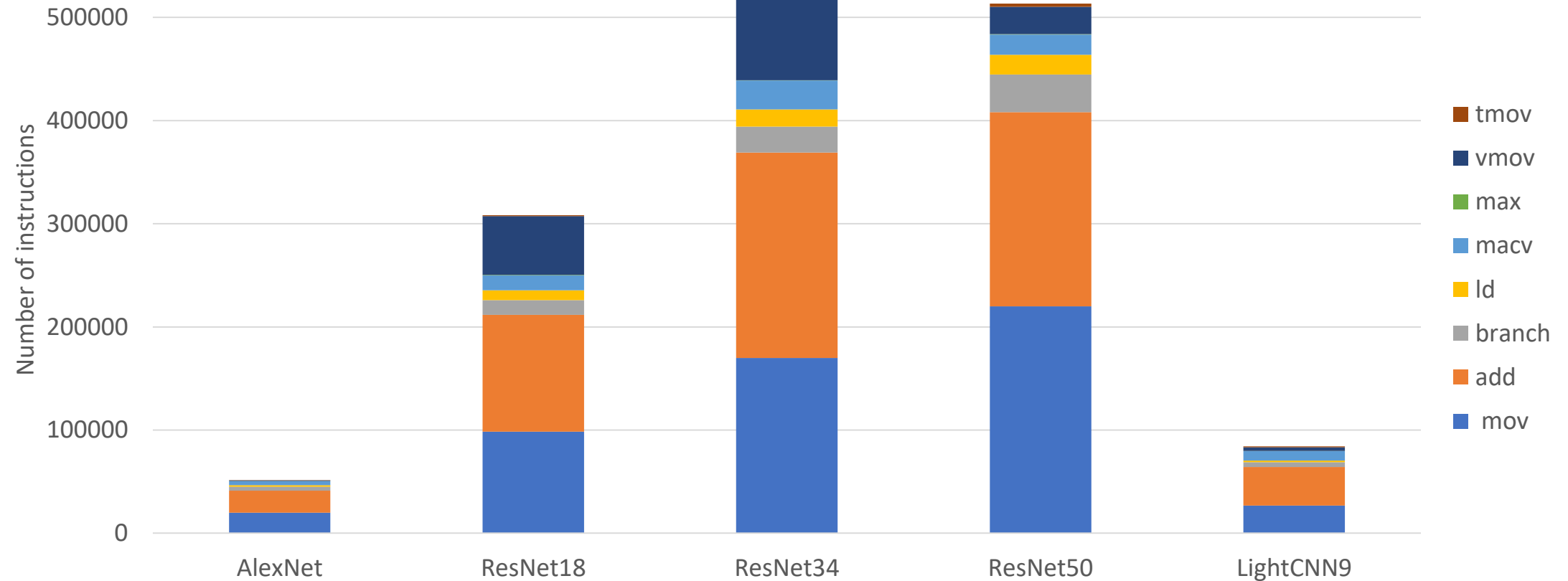
...



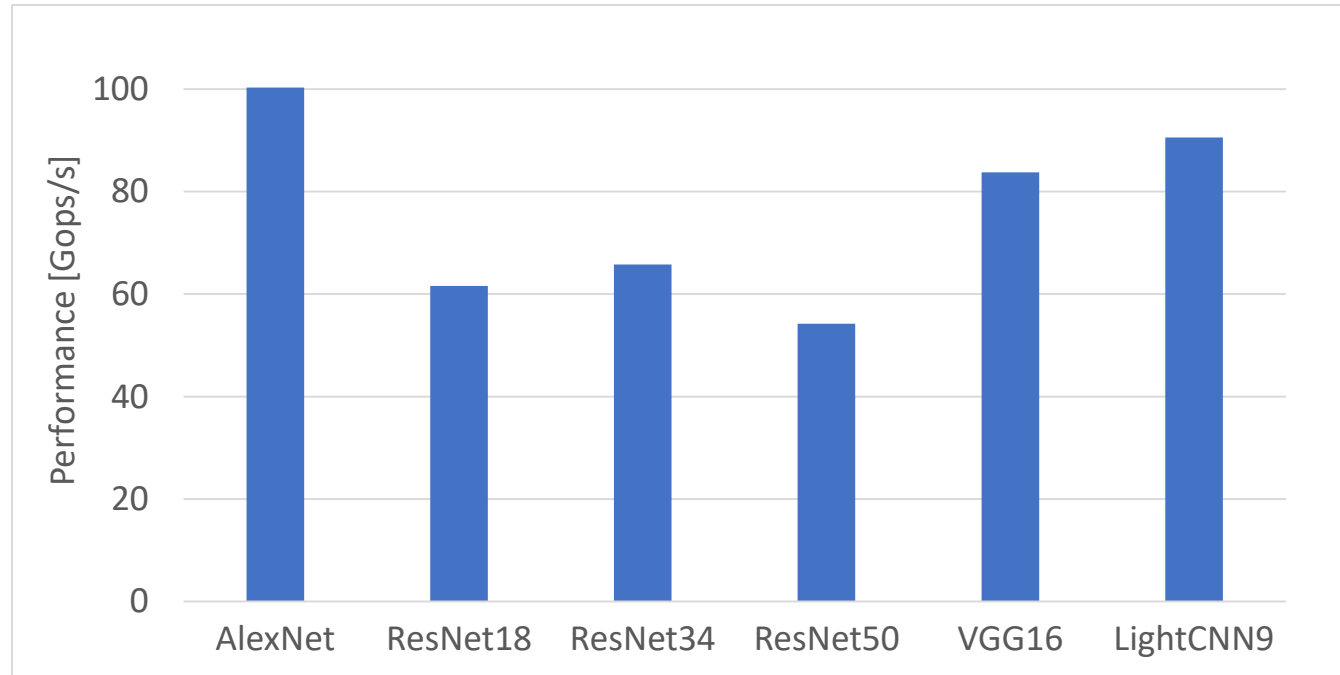
# Code generate



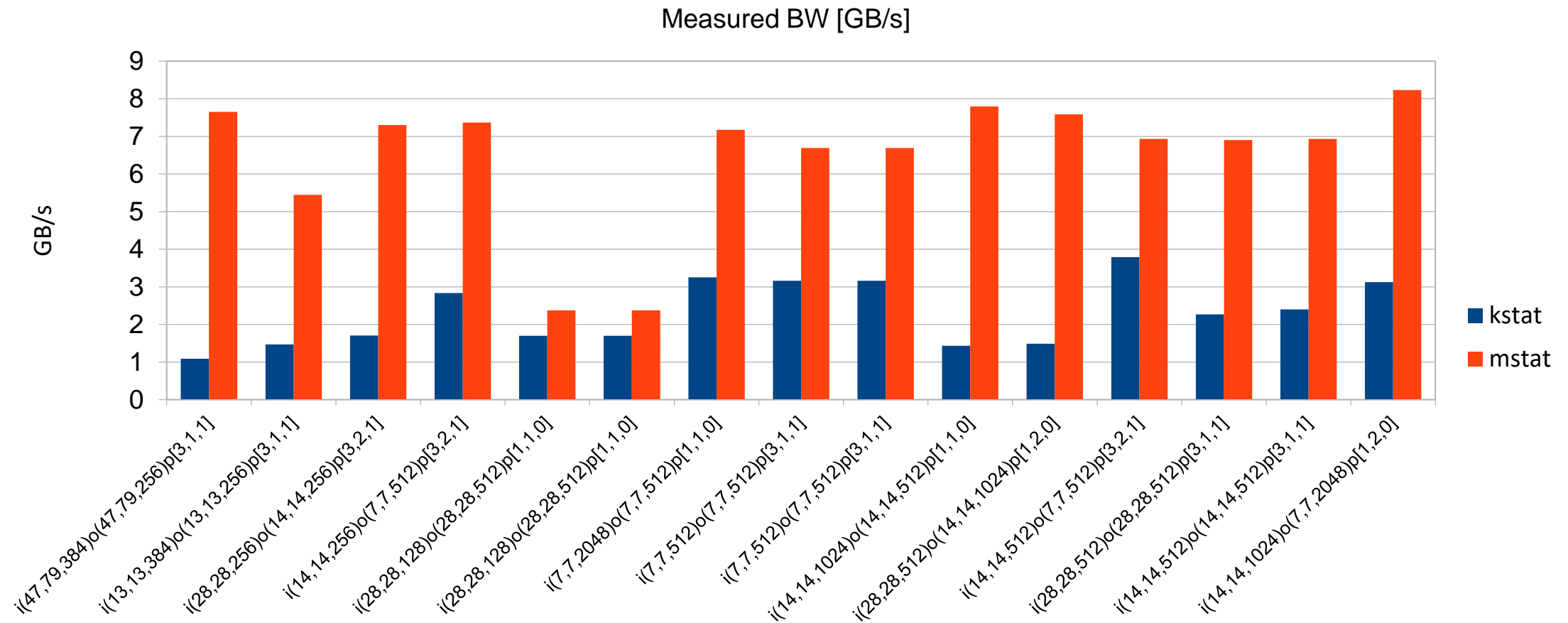
# Results – code



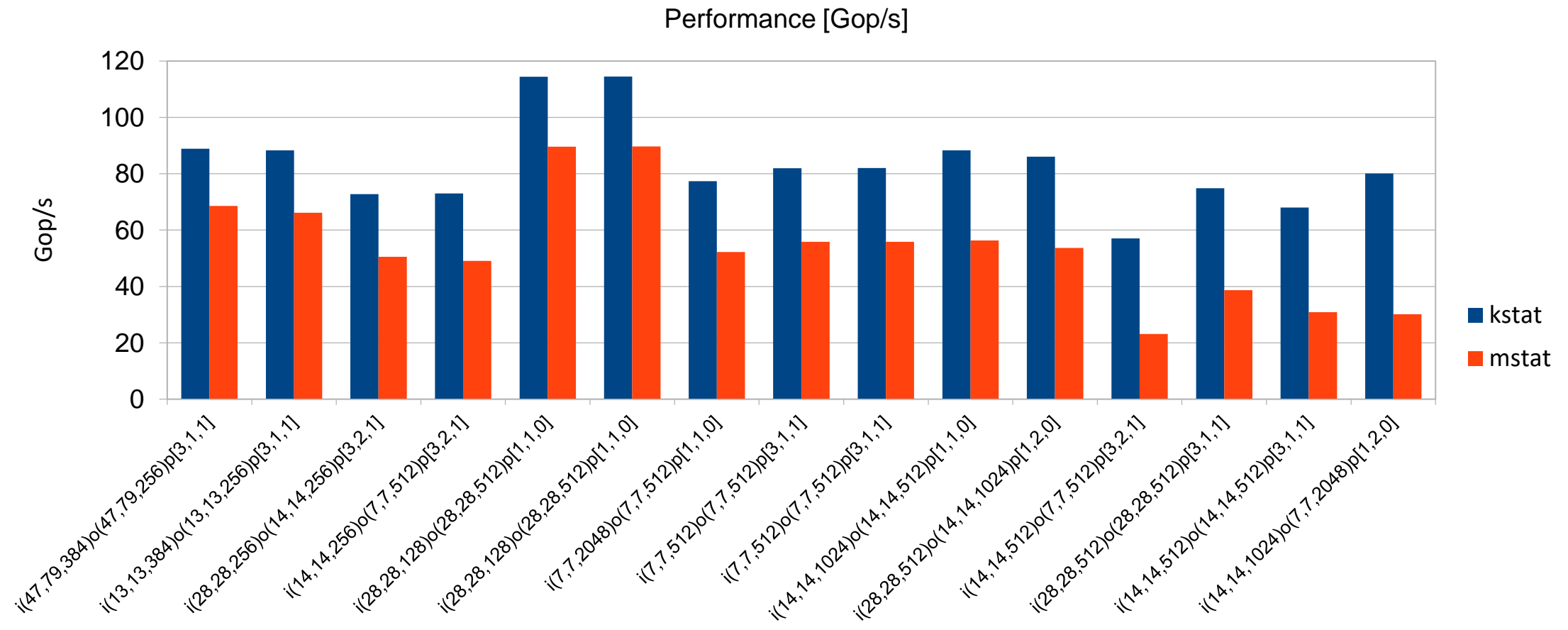
# Results - Performance



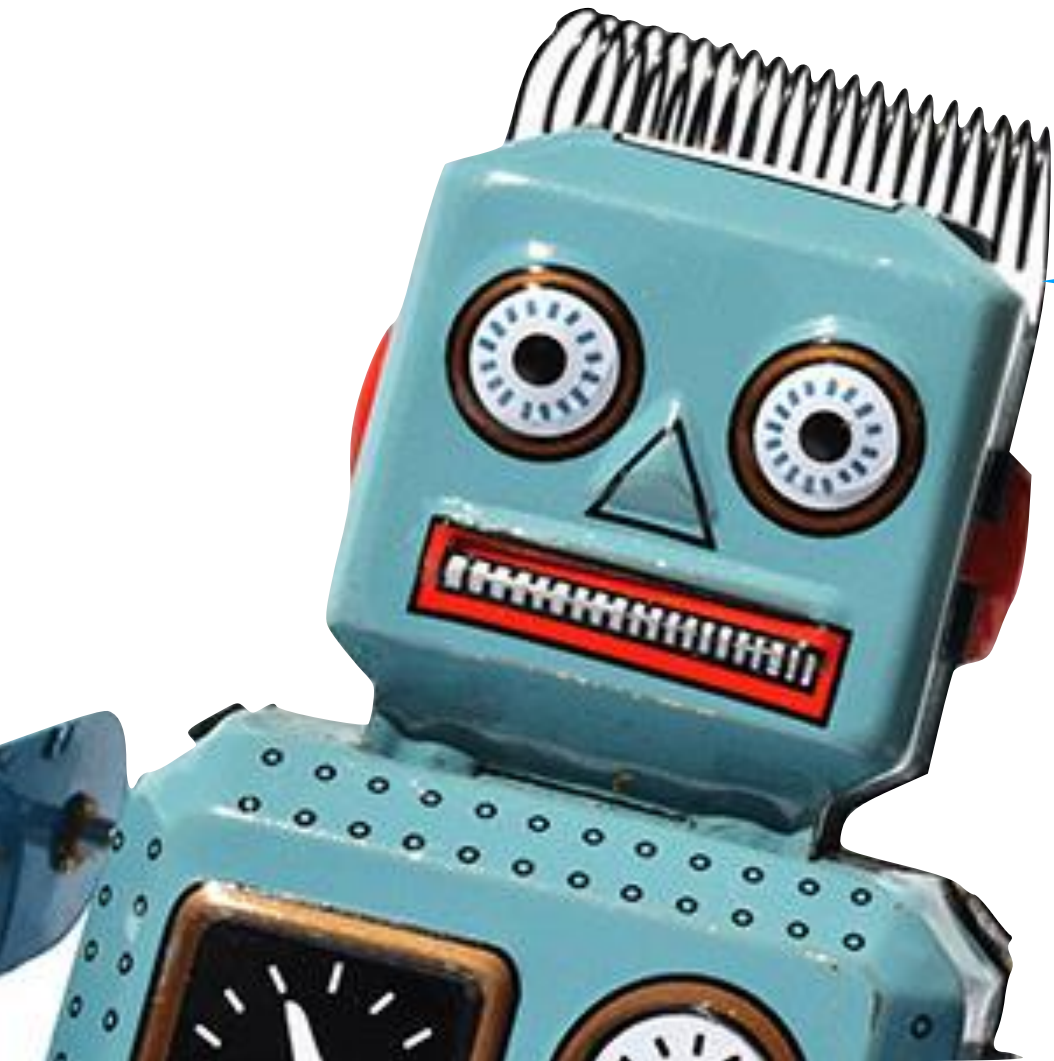
# Results Bandwidth



# Results







Thank you



FW  NXT