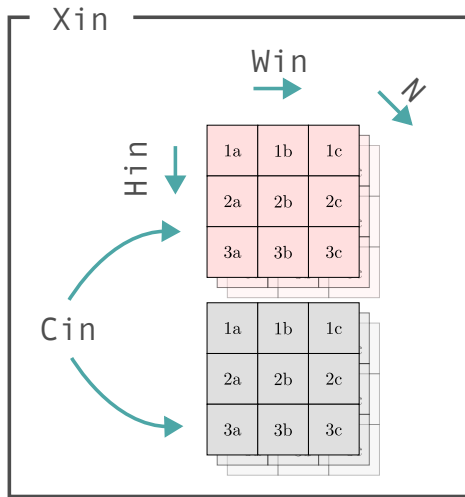# More convolutional nets

## Lecture 09 — CS 577 Deep Learning
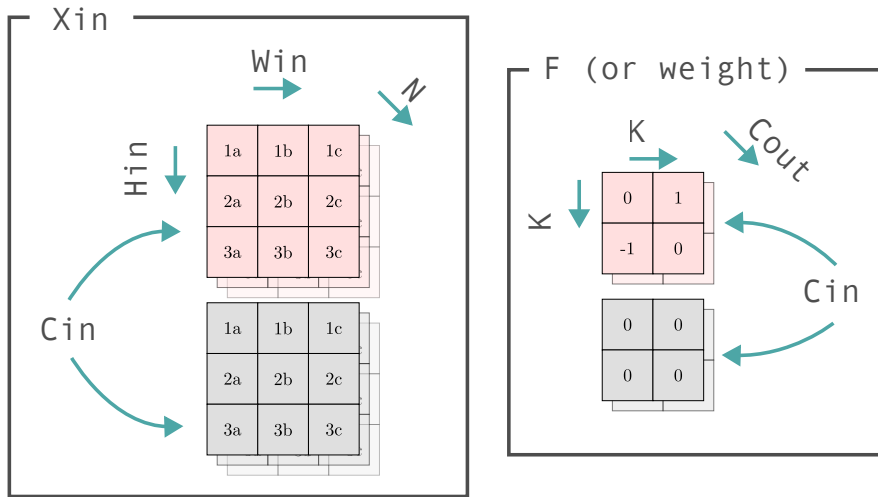
Instructor: Yutong Wang

Computer Science
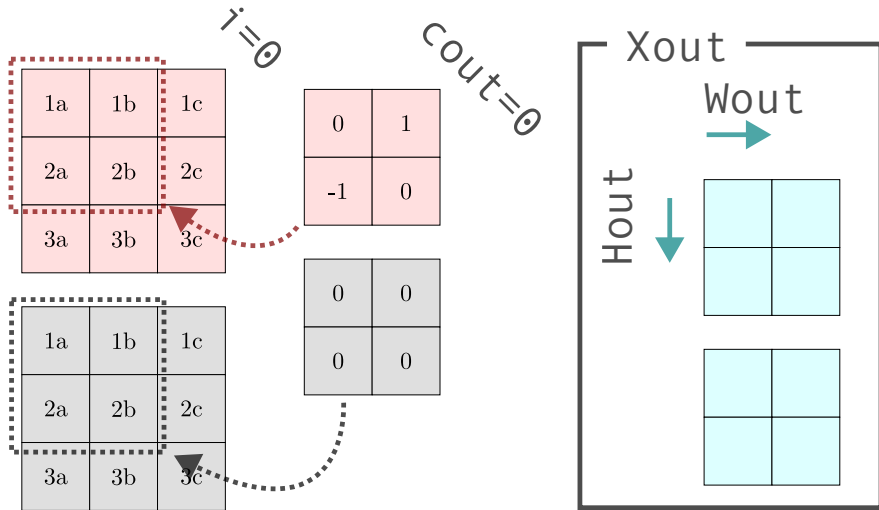Illinois Institute of Technology

October 16, 2024

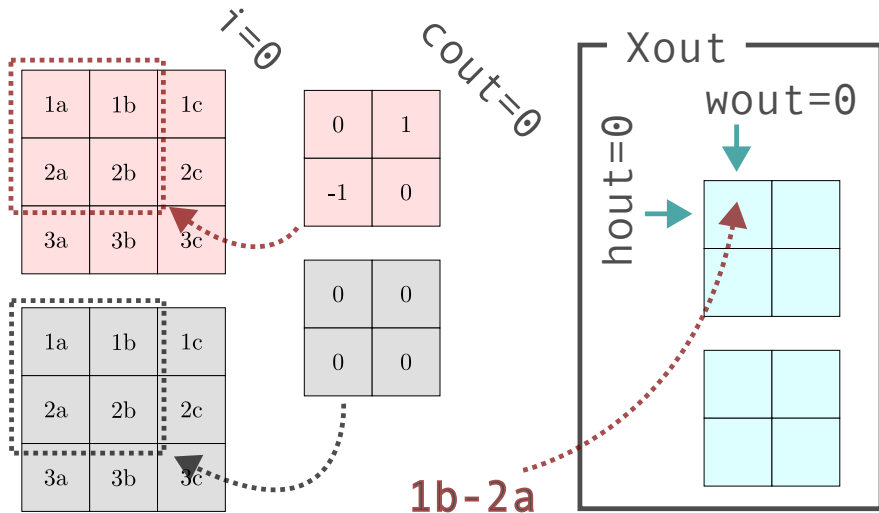# Input image tensor with shape $(N, C_{\text{in}}, H_{\text{in}}, W_{\text{in}})$

# A single convolution layer $(C_{\text{out}}, C_{\text{in}}, K, K)$

# Output: Xout with shape $(N, C_{\mathrm{out}}, H_{\mathrm{out}}, W_{\mathrm{out}})$

i=0

cout=0

Xout

wout=0

hout=0

| 0 | 1 |
| -1 | 0 |

| 0 | 0 |
| 0 | 0 |

| 1a | 1b | 1c |
| 2a | 2b | 2c |
| 3a | 3b | 3c |

| 1a | 1b | 1c |
| 2a | 2b | 2c |
| 3a | 3b | 3c |

**1b-2a**

wker=0

hker=0

| 0 | 1 |
| -1 | 0 |

| 0 | 0 |
| 0 | 0 |

Xout

wout=0

hout=0

**1b-2a**

patch_idx = 0

wker=1

hker=0

Xout

wout=0

hout=0

1b-2a

1b-2a

wker=1

hker=1

Xout

wout=0

hout=0

1b-2a

1c-2b

# patch_idx = 2



i=0

cout=0

Xout

wout=0

hout=1

| 1a | 1b | 1c |
| 2a | 2b | 2c |
| 3a | 3b | 3c |

| 0 | 1 |
| -1 | 0 |

| 0 | 0 |
| 0 | 0 |

| 1a | 1b | 1c |
| 2a | 2b | 2c |
| 3a | 3b | 3c |

2b-3a

# `patch_idx = 3`



i=0

cout=0

Xout

wout=1

hout=1

| 1a | 1b | 1c |
|----|----|----|
| 2a | 2b | 2c |
| 3a | 3b | 3c |

| 0 | 1 |
|---|---|
| -1 | 0 |

| 1a | 1b | 1c |
|----|----|----|
| 2a | 2b | 2c |
| 3a | 3b | 3c |

| 0 | 0 |
|---|---|
| 0 | 0 |

**2c-3b**

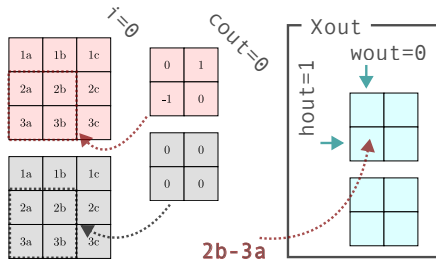# Sliding window (simple but slow)



```
1  # Suppose i, c_out, h_out, w_out are already defined, and
2  # X_out is initialized to all zeros
3  for cin in range(Cin):  # Input channels
4      for hker in range(K):  # Kernel height
5          for wker in range(K):  # Kernel width
6              Xout[i, cout, hout, wout] += (
7                  F[cout, cin, hker, wker] *
8                  Xin[i, cin, hout + hker, wout + wker])
```

There's gotta be a better way! "im2col" (Next)

# "im2col"

Main idea (in theory)

- Keep the window fixed in one spot
- Move the image tensor

Main idea (in theory)

- Create patches of the data
- Each patch has the same shape as the convolution filter
- Do matrix multiplication between the (flattened) patches and the (flattened) convolutional filter weights

# Flatten the weights

| 1a | 1b | 1c |
|----|----|----|
| 2a | 2b | 2c |
| 3a | 3b | 3c |

| 0 | 1 |
|----|----|
| -1 | 0 |

weight_flat

| 1a | 1b | 1c |
|----|----|----|
| 2a | 2b | 2c |
| 3a | 3b | 3c |

| 0 | 0 |
|----|----|
| 0 | 0 |

# Flatten the weights

# Flatten the weights

| 1a | 1b | 1c |
|----|----|----|
| 2a | 2b | 2c |
| 3a | 3b | 3c |

| 0  | 1 |
|----|---|
| -1 | 0 |

**weight_flat**

| 0 | 1 | -1 | 0 | 0 | 0 | 0 | 0 |
|---|---|----|---|---|---|---|---|

| 1a | 1b | 1c |
|----|----|----|
| 2a | 2b | 2c |
| 3a | 3b | 3c |

| 0 | 0 |
|---|---|
| 0 | 0 |

**Xin_patches_flat**

# Flatten the weights

# Flatten the weights

# Flatten the weights

# Flatten the weights



weight_flat

| 0 | 1 | -1 | 0 | 0 | 0 | 0 | 0 |
|---|---|----|---|---|---|---|---|

Xin_patches_flat

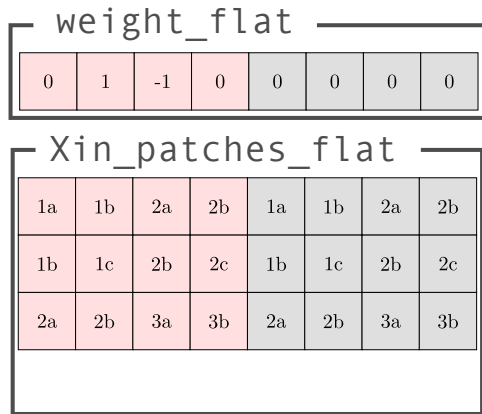| 1a | 1b | 2a | 2b | 1a | 1b | 2a | 2b |
|----|----|----|----|----|----|----|----|
| 1b | 1c | 2b | 2c | 1b | 1c | 2b | 2c |
| 2a | 2b | 3a | 3b | 2a | 2b | 3a | 3b |
| 2b | 2c | 3b | 3c | 2b | 2c | 3b | 3c |

# Sliding window via matmul

# Sliding window via matmul

# Flattening `Xin` into `Xin_flat`

# Flattening `Xin` into `Xin_flat`

Recall patch_idx=0:

| | | |
|---|---|---|
| 1a | 1b | 1c |
| 2a | 2b | 2c |
| 3a | 3b | 3c |

| | | |
|---|---|---|
| 1a | 1b | 1c |
| 2a | 2b | 2c |
| 3a | 3b | 3c |

Our goal:

— Xin_patches_flat —

| 1a | 1b | 2a | 2b | 1a | 1b | 2a | 2b |
|---|---|---|---|---|---|---|---|

What we have:

— Xin_flat —

| 1a | 1b | 1c | 2a | 2b | 2c | 3a | 3b | 3c | 1a | 1b | 1c | 2a | 2b | 2c | 3a | 3b | 3c |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

# The "im2col" matrix

# The "im2col" matrix



Xin_flat @

Our goal:

Xin_patches_flat

| 1a | 1b | 2a | 2b | 1a | 1b | 2a | 2b |

What we have:

Xin_flat

| 1a | 1b | 1c | 2a | 2b | 2c | 3a | 3b | 3c | 1a | 1b | 1c | 2a | 2b | 2c | 3a | 3b | 3c |

# The "im2col" matrix

Our goal:

Xin_patches_flat

| 1a | 1b | 2a | 2b | 1a | 1b | 2a | 2b |

Xin_flat @

What we have:

Xin_flat

| 1a | 1b | 1c | 2a | 2b | 2c | 3a | 3b | 3c | 1a | 1b | 1c | 2a | 2b | 2c | 3a | 3b | 3c |

# The "im2col" matrix



Our goal:

Xin_patches_flat

| 1a | 1b | 2a | 2b | 1a | 1b | 2a | 2b |

Xin_flat @

What we have:

Xin_flat

| 1a | 1b | 1c | 2a | 2b | 2c | 3a | 3b | 3c | 1a | 1b | 1c | 2a | 2b | 2c | 3a | 3b | 3c |

# The "im2col" matrix

# The "im2col" matrix

Recall patch_idx=1:

Our goal:

Xin_patches_flat

| 1b | 1c | 2b | 2c | 1b | 1c | 2b | 2c |
|----|----|----|----|----|----|----|----|

What we have:

Xin_flat

| 1a | 1b | 1c | 2a | 2b | 2c | 3a | 3b | 3c | 1a | 1b | 1c | 2a | 2b | 2c | 3a | 3b | 3c |

# The "im2col" matrix

Xin_flat @

Our goal:

Xin_patches_flat

| 1b | 1c | 2b | 2c | 1b | 1c | 2b | 2c |

What we have:

Xin_flat

| 1a | 1b | 1c | 2a | 2b | 2c | 3a | 3b | 3c | 1a | 1b | 1c | 2a | 2b | 2c | 3a | 3b | 3c |

# The "im2col" matrix

Our goal:

Xin_patches_flat

| 2a | 2b | 3a | 3b | 2a | 2b | 3a | 3b |

Xin_flat @

What we have:

Xin_flat

| 1a | 1b | 1c | 2a | 2b | 2c | 3a | 3b | 3c | 1a | 1b | 1c | 2a | 2b | 2c | 3a | 3b | 3c |

# The "im2col" matrix

Xin_flat @

im2col

=

Xin_im2col

| 1a | 1b | 2a | 2b | 1a | 1b | 2a | 2b |

# The "im2col" matrix



Xin_flat @

im2col

=

Xin_im2col

| 1a | 1b | 2a | 2b | 1a | 1b | 2a | 2b | 1b | 1c | 2b | 2c | 1b | 1c | 2b | 2c |

# The "im2col" matrix



Xin_flat @

im2col

=

Xin_im2col

# The "im2col" matrix

Xin_flat @

im2col

=

Xin_im2col

# The "im2col" matrix



Xin_flat @

im2col

=

Xin_im2col

# The "im2col" matrix



Xin

Xin_flat

Xin_flat @

| 1a | 1b | 1c | 2a | 2b | 2c | 3a | 3b | 3c | 1a | 1b | 1c | 2a | 2b | 2c | 3a | 3b | 3c |

im2col

=

Xin_im2col

output_index =
    patch_idx * patch_size
    + patch_position

# The "im2col" matrix

Xin_flat @

im2col

=

Xin_im2col

Xin

Xin_flat

patch_size =
    Cin * K * K

output_index =
    patch_idx * patch_size
    + patch_position

# The "im2col" matrix



Xin

Xin_flat @

Xin_flat

| 1a | 1b | 1c | 2a | 2b | 2c | 3a | 3b | 3c | 1a | 1b | 1c | 2a | 2b | 2c | 3a | 3b | 3c |

im2col

```
patch_position =
    cin * K * K
    + hker * K + wker
```

=

Xin_im2col

```
output_index =
    patch_idx * patch_size
    + patch_position
```

# The "im2col" matrix



input_index

Xin_flat @

im2col

=

Xin_im2col

Xin

Xin_flat

patch_position =
   cin * K * K
   + hker * K + wker

output_index =
   patch_idx * patch_size
   + patch_position

# The "im2col" matrix

input_index

input_index =
    cin * Hin * Win
    + hout * Win + wout
    + hker * Win + wker

Xin

Xin_flat @

Xin_flat

im2col

patch_position =
    cin * K * K
    + hker * K + wker

=

output_index =
    patch_idx * patch_size
    + patch_position

Xin_im2col

# The "im2col" matrix

input_index

input_index =
  cin * Hin * Win
  + hout * Win + wout
  + hker * Win + wker

Xin

Xin_flat @

Xin_flat

im2col

patch_position =
  cin * K * K
  + hker * K + wker

=

output_index =
  patch_idx * patch_size
  + patch_position

Xin_im2col

# The "im2col" matrix



input_index =
  cin * Hin * Win
  + hout * Win + wout
  + hker * Win + wker

input_index

Xin_flat @

im2col

=

Xin_im2col

Xin

Xin_flat

patch_position =
  cin * K * K
  + hker * K + wker

output_index =
  patch_idx * patch_size
  + patch_position

# The "im2col" matrix



input_index

Xin_flat @

im2col

=

Xin_im2col

input_index =
    cin * Hin * Win
    + hout * Win + wout
    + hker * Win + wker

Xin

Xin_flat

patch_position =
    cin * K * K
    + hker * K + wker

output_index =
    patch_idx * patch_size
    + patch_position

# The "im2col" matrix



input_index

Xin_flat @

im2col

=

Xin_im2col

input_index =
    cin * Hin * Win
    + hout * Win + wout
    + hker * Win + wker

Xin

Xin_flat

patch_position =
    cin * K * K
    + hker * K + wker

output_index =
    patch_idx * patch_size
    + patch_position

# The "im2col" matrix



input_index

input_index =
cin * Hin * Win
+ hout * Win + wout
+ hker * Win + wker

Xin_flat @

im2col

Xin

Xin_flat

patch_position =
    cin * K * K
    + hker * K + wker

=

output_index =
    patch_idx * patch_size
    + patch_position

Xin_im2col

# The "im2col" matrix



input_index =
cin * Hin * Win
+ hout * Win + wout
+ hker * Win + wker

Xin

Xin_flat @

Xin_flat

im2col

patch_position =
cin * K * K
+ hker * K + wker

=

output_index =
patch_idx * patch_size
+ patch_position

Xin_im2col

# Putting it all together

```python
1  def im2col_matrix_dense(Xin, K, S=1):
2      N, Cin, Hin, Win = Xin.shape
3      Hout, Wout = Hin - K + 1, Win - K + 1
4      P = Hout * Wout   # Total number of patches per image
5      patch_size = Cin * K * K # Size of each flattened patch
6      im2col_mat_dense = np.zeros((Cin * Hin * Win, P * patch_size))
7
8      # [main loop on next slide...]
9
10     return im2col_mat_dense
```
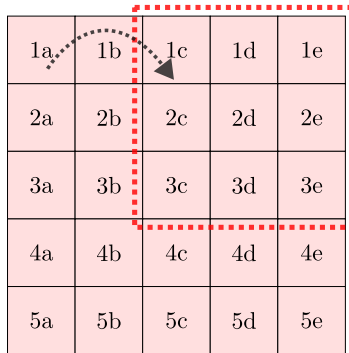
# Putting it all together

```
1    # [ continued from previous slide ...]
2    patch_idx = 0
3    for hout in range ( Hout ):
4        for wout in range ( Wout ):
5            for cin in range ( Cin ):
6                for hker in range ( K ):
7                    for wker in range ( K ):
8                        input_index = cin * Hin * Win + hout * Win +
     wout + hker * Win + wker
9                        patch_position = cin * K * K + hker * K + wker
10                       output_index = patch_idx * patch_size +
     patch_position
11                       im2col_mat_dense [ input_index , output_index ] = 1
12           patch_idx += 1
```

| 1a | 1b | 1c | 1d | 1e |
|----|----|----|----|----|
| 2a | 2b | 2c | 2d | 2e |
| 3a | 3b | 3c | 3d | 3e |
| 4a | 4b | 4c | 4d | 4e |
| 5a | 5b | 5c | 5d | 5e |

| 1a | 1b | 1c | 1d | 1e |
|----|----|----|----|----|
| 2a | 2b | 2c | 2d | 2e |

# Stride = 2

# Implement `im2col_matrix_dense` with stride

See `lec09-in-class-ex1-im2col.ipynb`
End of part 1.

# Goal for today: implement LeNet5



Image from LeCun et al. 1998

# `im2col.shape` = $C_{\text{in}} \cdot H_{\text{in}} \cdot W_{\text{in}}$**-by-**$P \cdot C_{\text{out}} \cdot K^2$



Xin_flat @

im2col

=

Xin_im2col

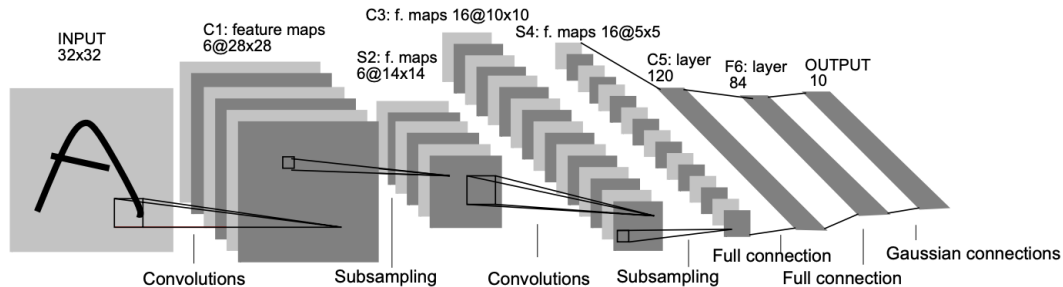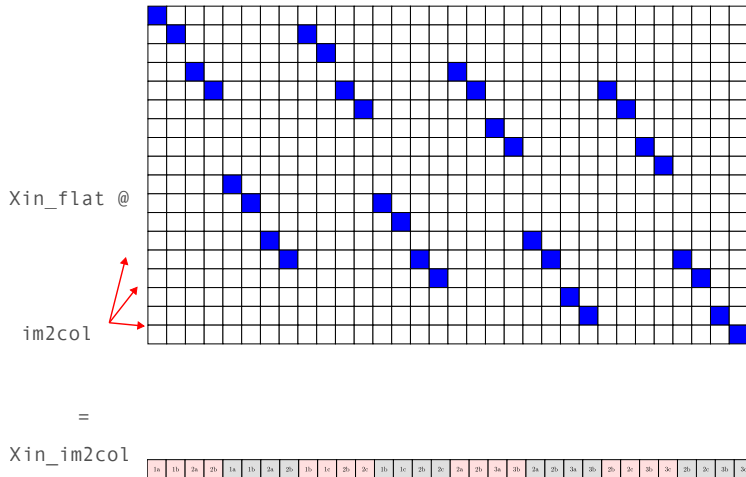# `im2col.shape` = $C_{\text{in}} \cdot H_{\text{in}} \cdot W_{\text{in}}$**-by-**$P \cdot C_{\text{in}} \cdot K^2$

For 1st convolution layer in LeNet5.

- $H_{\text{in}} = W_{\text{in}} = 32$
- $C_{\text{in}} = 1$
- $K = 5$
- $S = 1$
- $H_{\text{out}} = H_{\text{in}} - K + 1 = 28$
- $W_{\text{out}} = W_{\text{in}} - K + 1 = 28$
- $P = 28 * 28$
- Final shape of $\text{im2col} = 1024$ -by- $19600$
- $20{,}070{,}400$ float64's to store $\approx 160$ MB
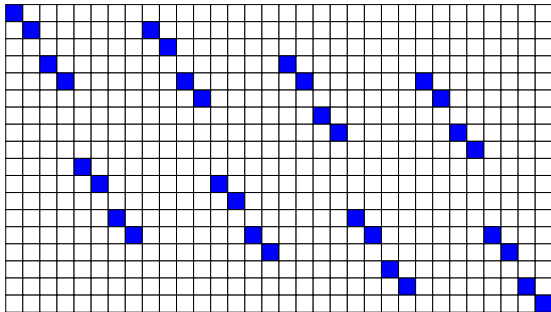
# Sparse matrices

Pros

- Store just the non-zero values of the matrix (memory efficient)
- Faster than dense matmul when matrices are mostly zeros (speedier)

Cons

- Irregular memory access (less cache-optimization friendly)
- Nightmare for GPUs (as of now, to the best of my knowledge)

# from scipy.sparse import csr_matrix



```
1    data = [1,1,1,1,1,...] # len(data) = P * patch_size
2    row_indices = [0,1,3,4,9,...]
3    col_indices = [0,1,2,3,4,...]
4    im2col_mat_sparse = csr_matrix((data,
5                                    (row_indices, col_indices)),
6                                   shape=(n_rows, n_cols))
```

# Example

```
1 data = [1,2,3,1]
2 row_indices = [1,2,0,1]
3 col_indices = [0,1,2,1]
4 sparse_mat_example = csr_matrix((data, (row_indices, col_indices)),
    shape=(3, 3))
5 sparse_mat_example.toarray()
```

Output:

```
1 array([[0, 0, 3],
2        [1, 1, 0],
3        [0, 2, 0]])
```

# Dense-sparse multiplication

```
1 X = np.arange(9).reshape(3,3)
2 ## this won't work
3 # np.matmul(X, sparse_mat_example)
4 ## this works bc method resolution order
5 X                 @     sparse_mat_example
```

```
1 array([[0, 1, 2],      array([[0, 0, 3],      array([[ 1,  5,  0],
2        [3, 4, 5],   @         [1, 1, 0],   =         [ 4, 14,  9],
3        [6, 7, 8]])            [0, 2, 0]])            [ 7, 23, 18]])
```

# Implement `im2col_matrix_sparse`

- `lec09-in-class-ex2-sparse.ipynb`

# Quick update

Can add/mul constants both on left and right now.

```python
def __radd__(self, other):
    return self + other

def __rmul__(self, other):
    return self * other

def __rsub__(self, other):
    return (-self) + other

def __rtruediv__(self, other):
    return ag.Tensor(other) / self
```

# Quick update

```
1    def moveaxis(input, source, destination):
2        output = ag.Tensor(np.moveaxis(input.value, source, destination)
     , inputs=[input], op="moveaxis")
3
4        def _backward():
5            input.grad += np.moveaxis(output.grad, source, destination)
6            return None
7        output._backward = _backward
8        return output
```
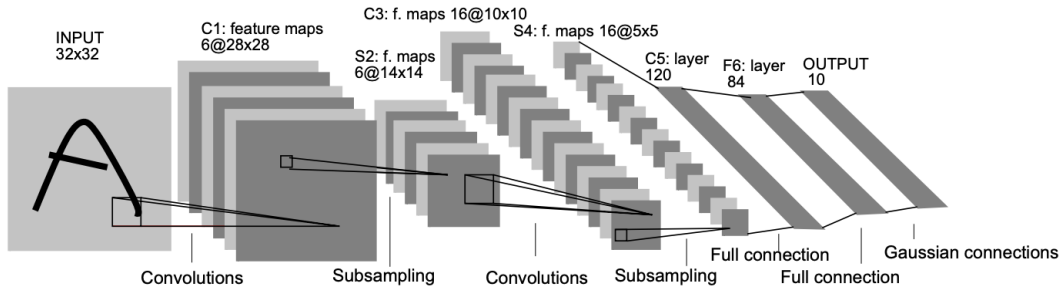
# Enabling autograd for dense-sparse matmul

- Problem: our current autograd-enabled matmul only supports when both inputs are numpy arrays
- Observation: we don't need to track gradients on the im2col_mat
- Solution:

```
 1    # [sp]arse [c]onstant (non-AG-enabled) [mat]rix [mul]tiplication
 2    def spcmatmul(input, sparse_mat: csr_matrix):
 3        output = ag.Tensor(input.value @ sparse_mat,
 4                           inputs = [input],
 5                           op="spcmatmul")
 6
 7        def _backward():
 8            input.grad += output.grad @ sparse_mat.T
 9            return None
10
11        output._backward = _backward
12        return output
```
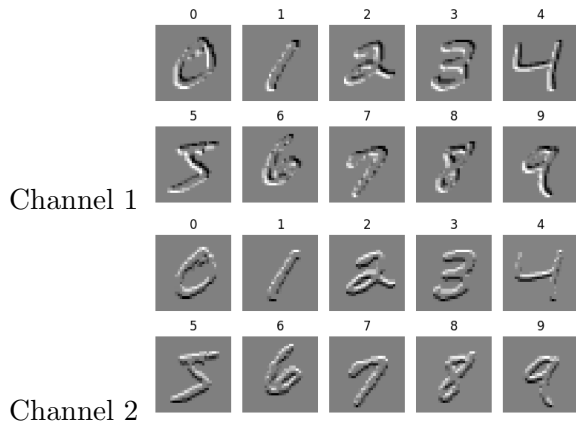
# Exercise 3

- Create the `Conv2d` layer
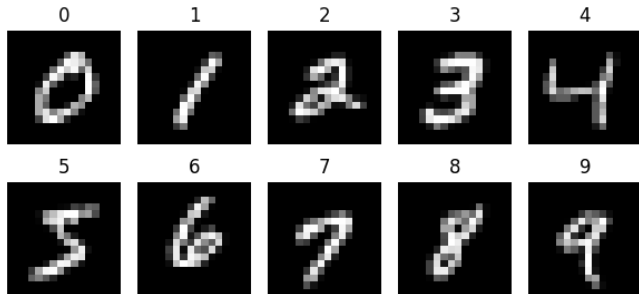- Create the `AvgPool2d` layer
- Create the `LeNet5` model



- Create the `nn.CrossEntropyLoss`

# Sanity check for conv2d



Channel 1

Channel 2

# Sanity check for pool

# Sanity check for training run

You should achieve 100 percentage training accuracy

# References I

[LeC+98]  Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner.
          "Gradient-based learning applied to document recognition". In:
          *Proceedings of the IEEE* 86.11 (1998), pp. 2278–2324.