

The CoRa Tensor Compiler: Compilation for Ragged Tensors With Minimal Padding

Pratik Fegade¹,
Tianqi Chen^{1,2}, Phillip B. Gibbons¹, Todd C. Mowry¹

¹Carnegie Mellon University

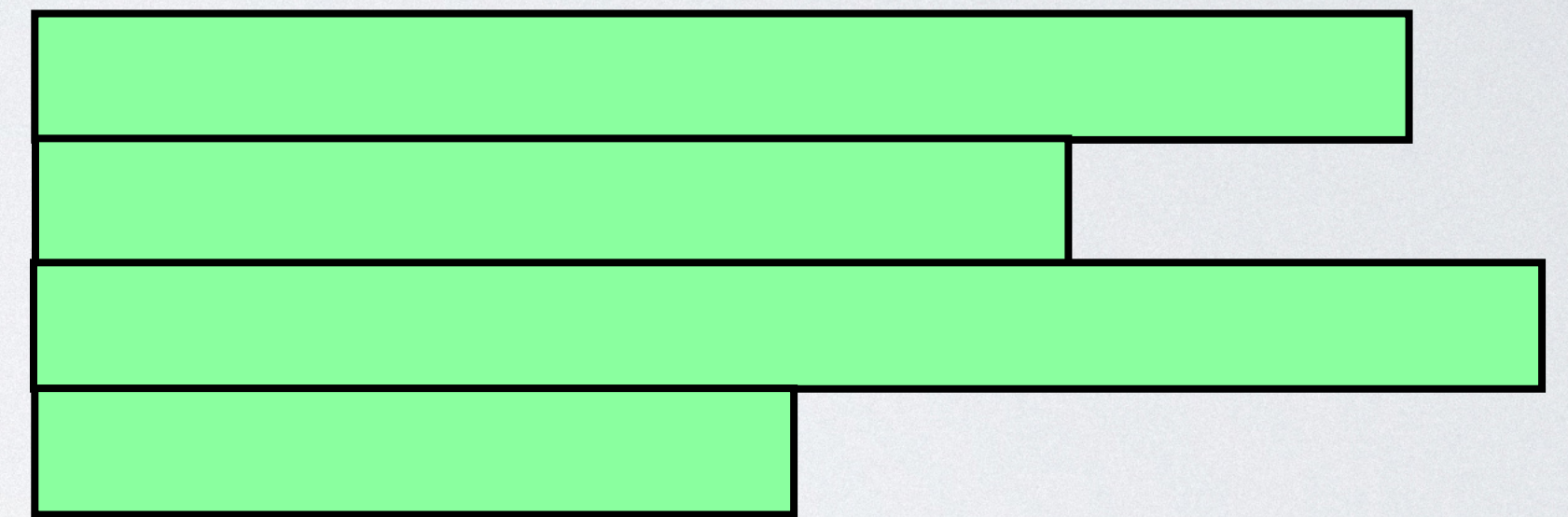
²OctoML

Ragged Tensors in Deep Learning

- Natural language processing

```
input_batch = [  
3  [Dogs, bark, .],  
5  [Maine, is, a, state, .],  
4  [The, song, rocks, !],  
1  [Hello]  
]
```

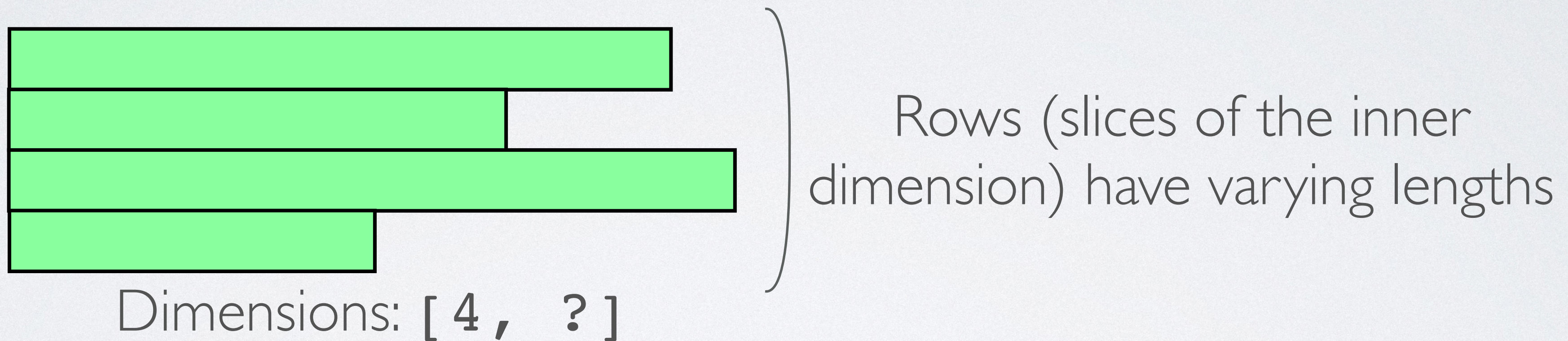
- Image processing



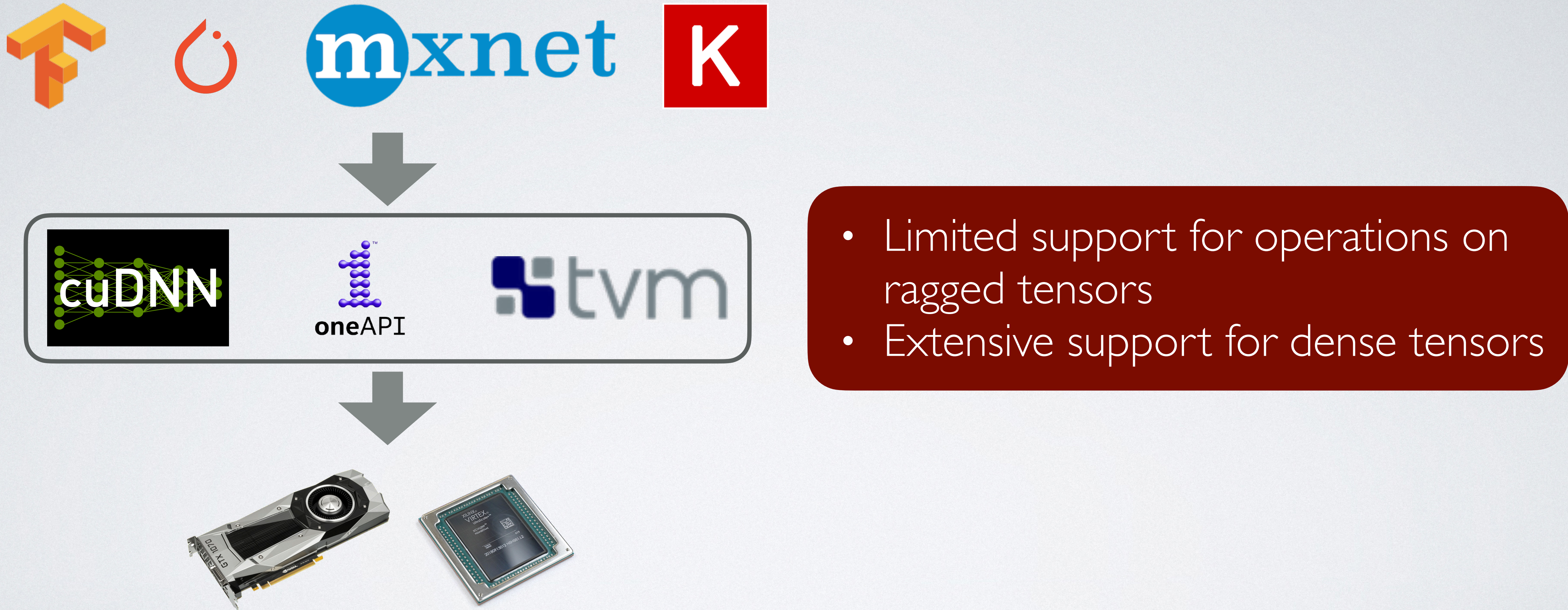
Ragged Tensor

Ragged Tensors

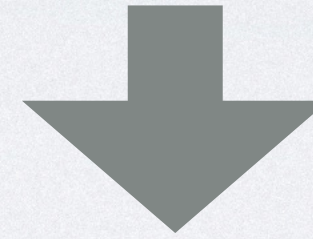
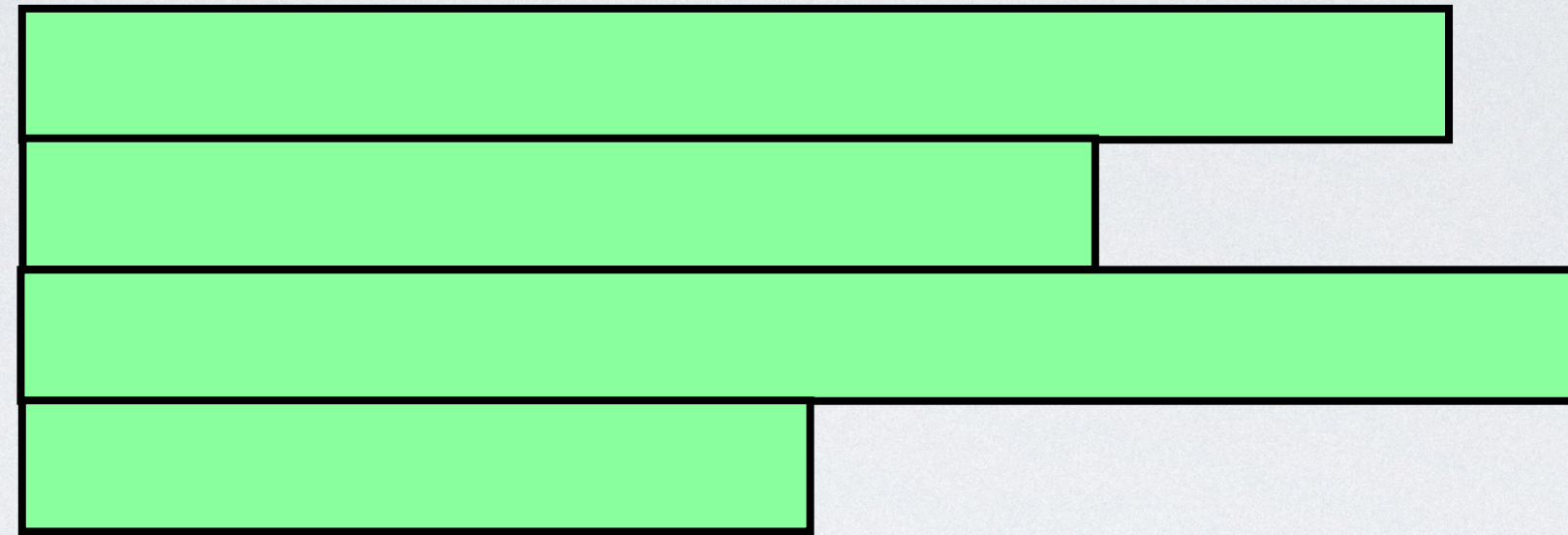
- Ragged tensor is a tensor where the slices corresponding to one or more dimensions have varying lengths



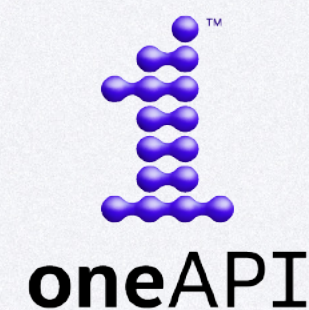
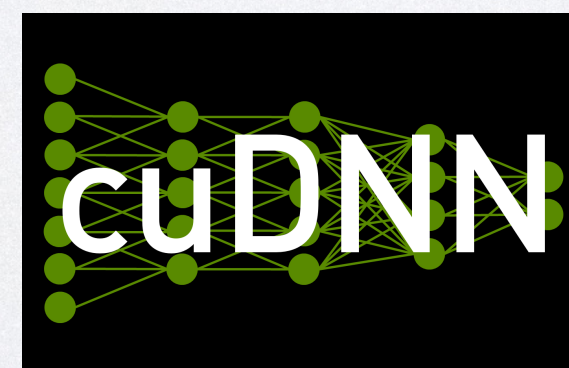
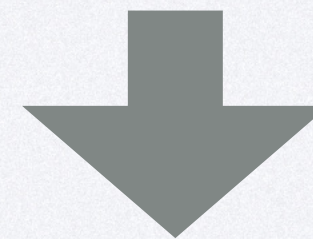
Limited Support for Ragged Tensor Operators



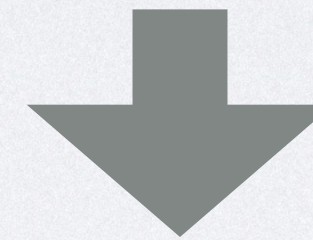
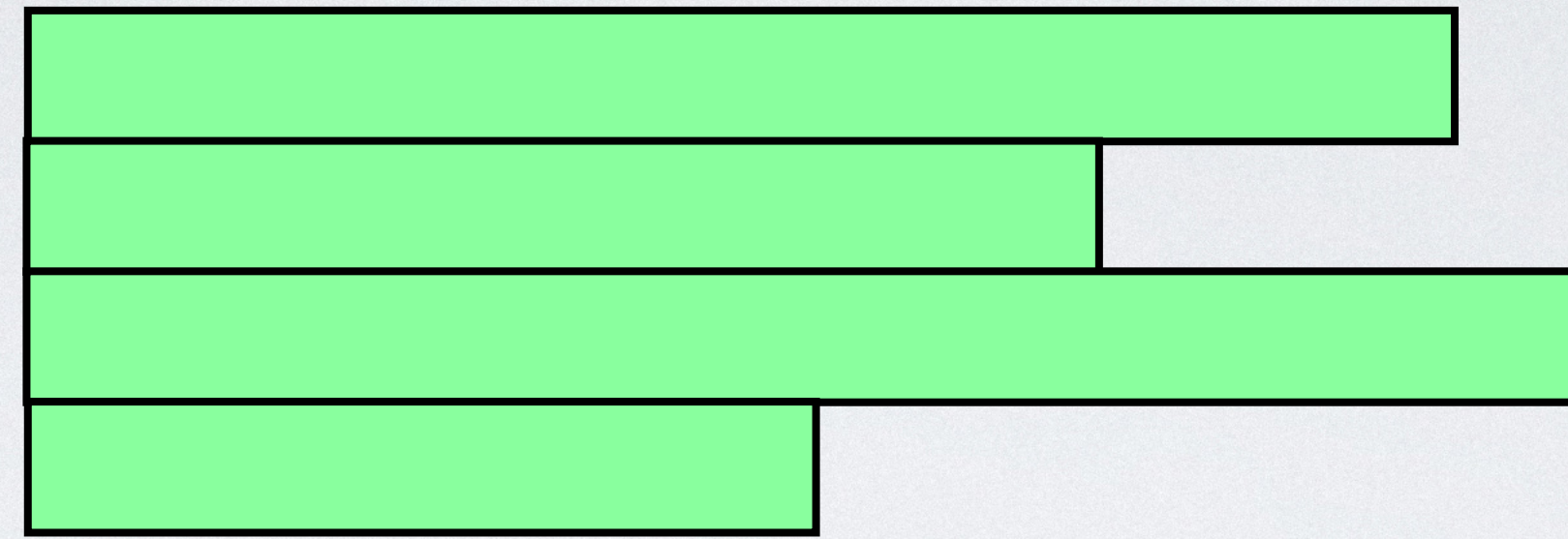
And Padding Leads to Wasted Computation



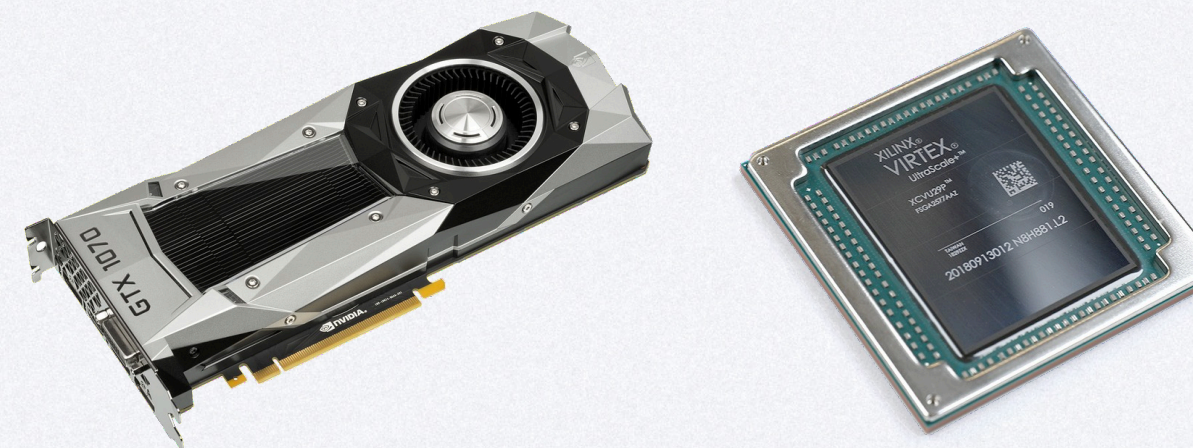
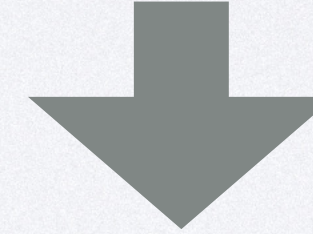
1.07 - 2.41X wasted computation!



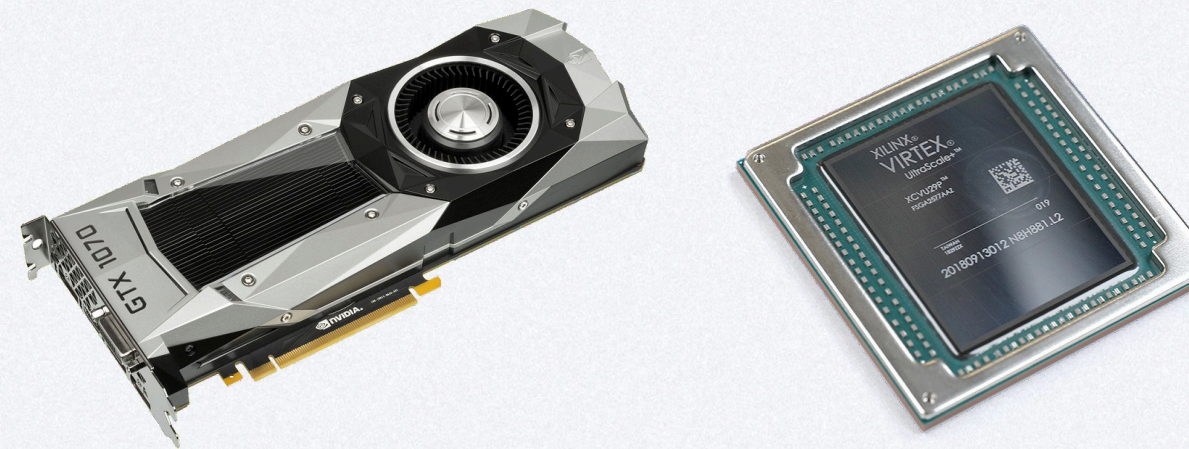
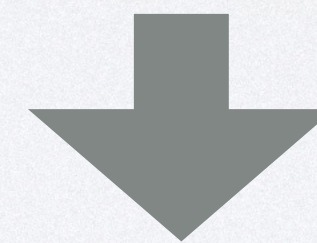
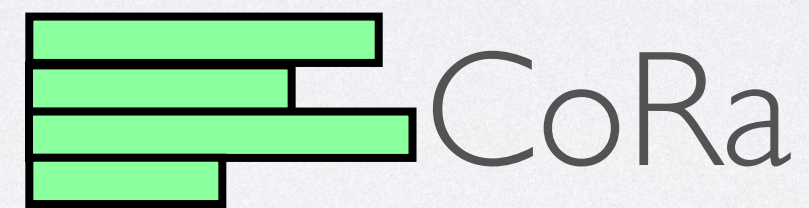
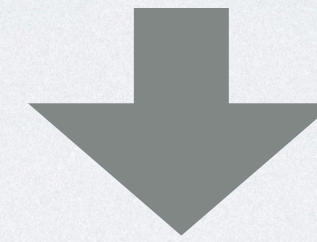
Ideal Execution: Compilation Without Padding



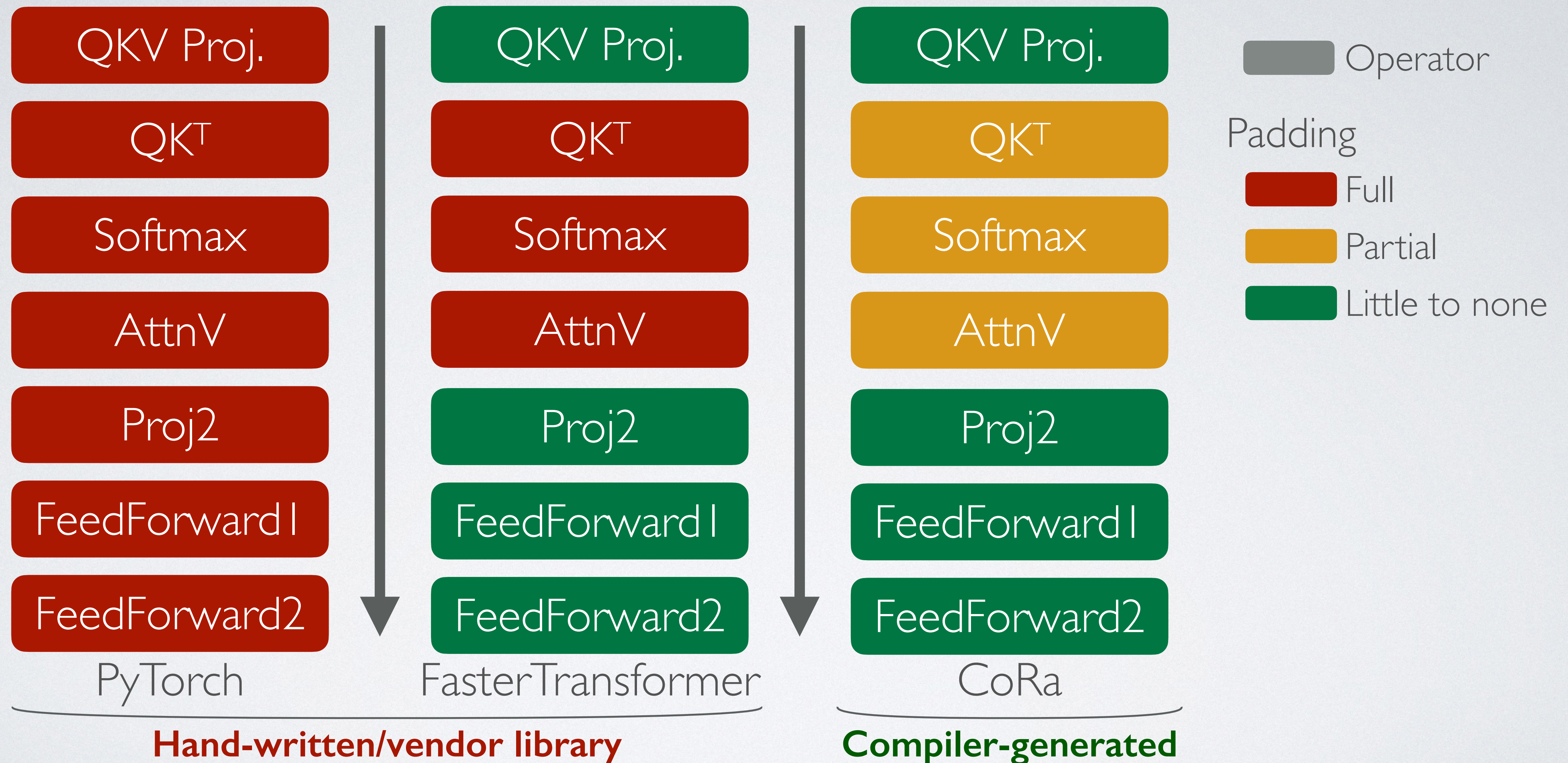
?



CoRa Enables Ragged Tensor Execution for Higher Frameworks



CoRa Enables Transformer Implementation Without Padding



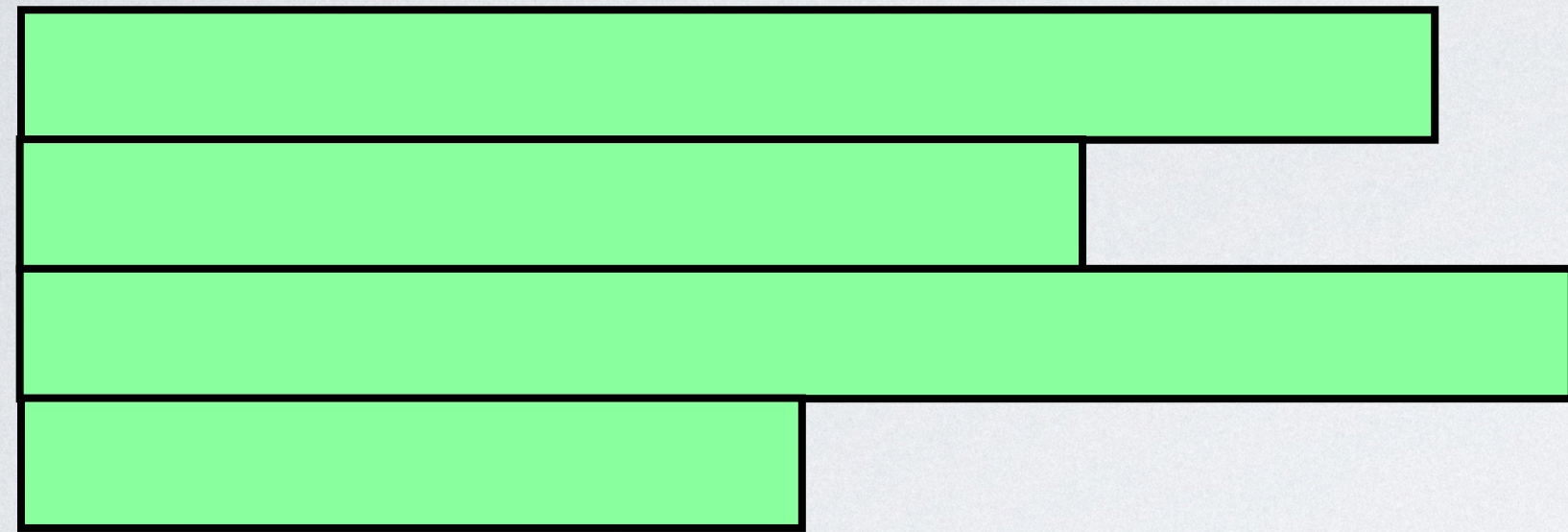
Outline

- Motivation: Inefficient Support for Ragged Tensors
- CoRa: Our Compiler Based Solution
 - Scheduling and lowering
 - API and overview
- Evaluation
- Wrapping up

Outline

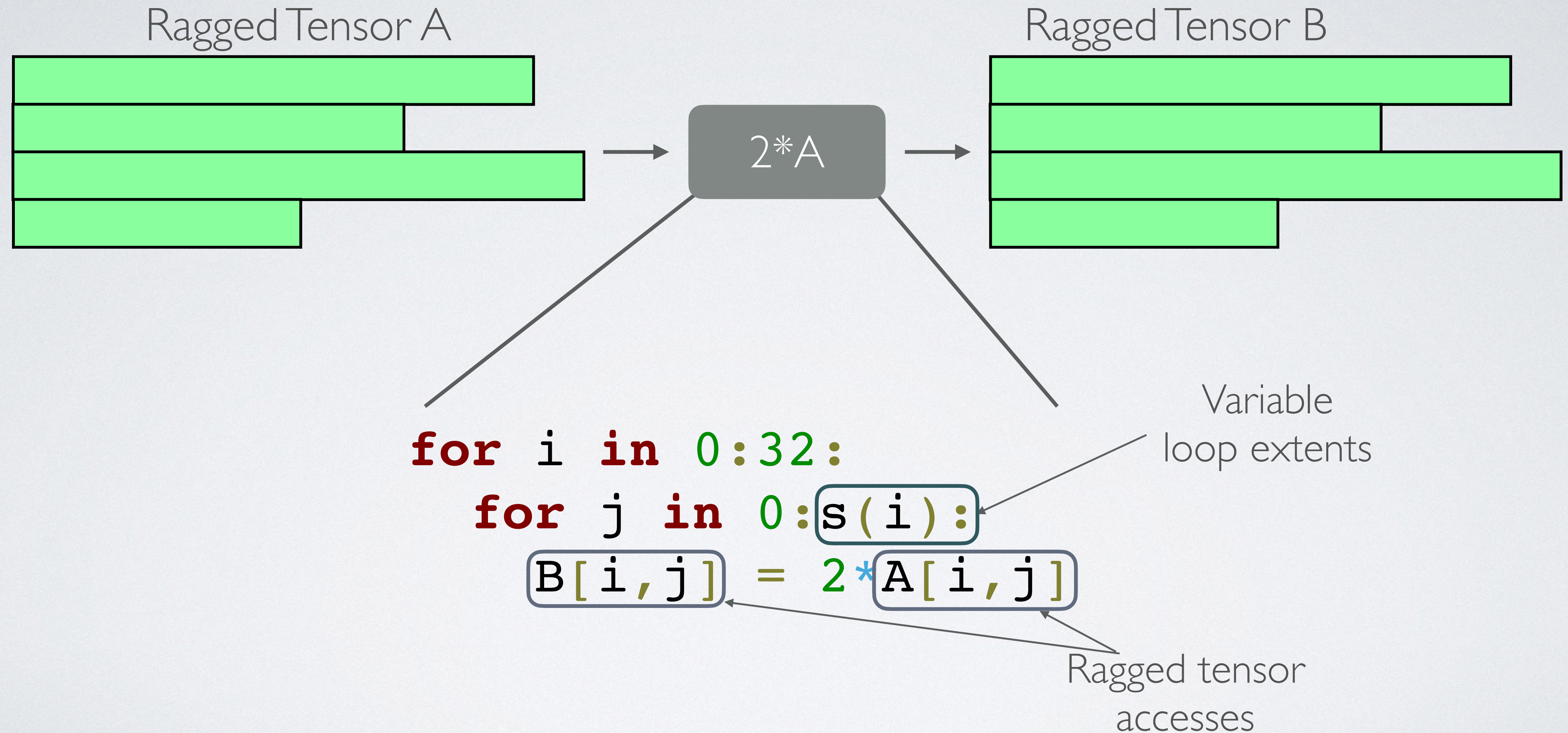
- Motivation: Inefficient Support for Ragged Tensors
- **CoRa: Our Compiler Based Solution**
 - Scheduling and lowering
 - API and overview
- Evaluation
- Wrapping up

Ragged Computations Are Similar to Dense Computations



Densely packed data with no holes, like dense tensors

Ragged Computations Are Similar to Dense Computations



Ragged Computations Are Similar to Dense Computations

- Densely packed data with no holes, like dense tensors
- Ragged computations are similar to dense tensor computations

Reuse abstractions and techniques from dense tensor compilers

```
for i in 0:32:  
  for j in 0:s(i):  
    B[i,j] = 2*A[i,j]
```

Variable
loop extents

Ragged tensor
accesses

Generalize

- Compiler's loop representations
- Scheduling primitives and their impl.

Generalize

- Tensor storage scheme
- Tensor access lowering

Ragged Computations Are Similar to Dense Computations

- Densely packed data with no holes, like dense tensors
- Ragged computations are similar to dense tensor computations

Reuse abstractions and techniques from dense tensor compilers

```
for i in 0:32:  
  for j in 0:s(i):  
    B[i,j] = 2*A[i,j]
```

Variable
loop extents

Ragged tensor
accesses

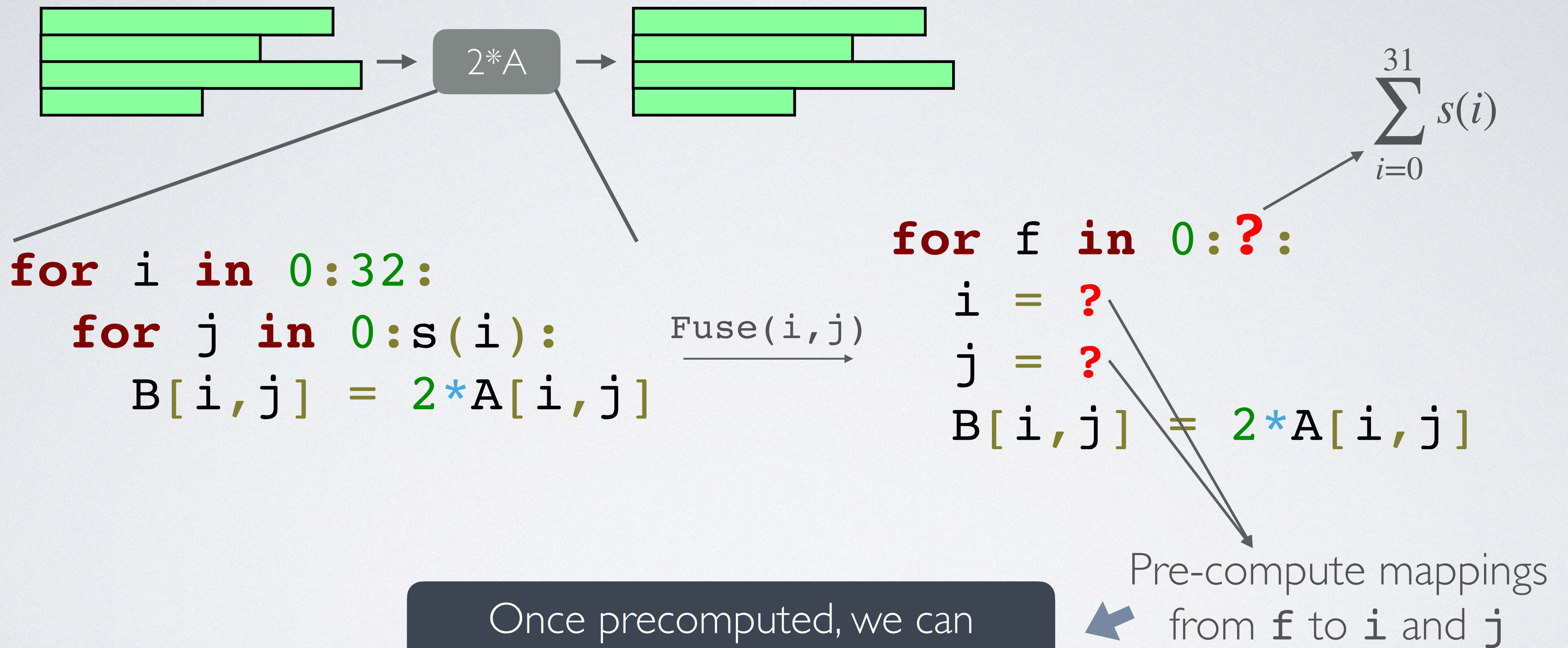
Generalize

- Compiler's loop representations
- Scheduling primitives and their impl.

Generalize

- Tensor storage scheme
- Tensor access lowering

Loop Fusion in Ragged Operators



Ragged Computations Are Similar to Dense Computations

- Densely packed data with no holes, like dense tensors
- Ragged computations are similar to dense tensor computations

Reuse abstractions and techniques from dense tensor compilers

```
for i in 0:32:
```

```
  for j in 0:s(i):
```

```
    B[i,j] = 2*A[i,j]
```

Variable
loop extents

Ragged tensor
accesses

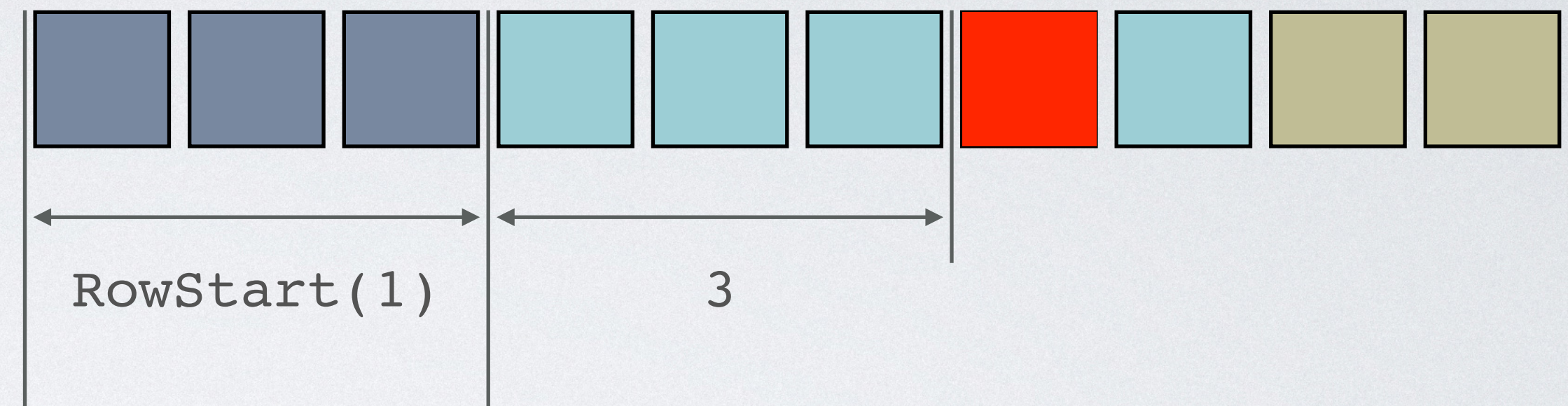
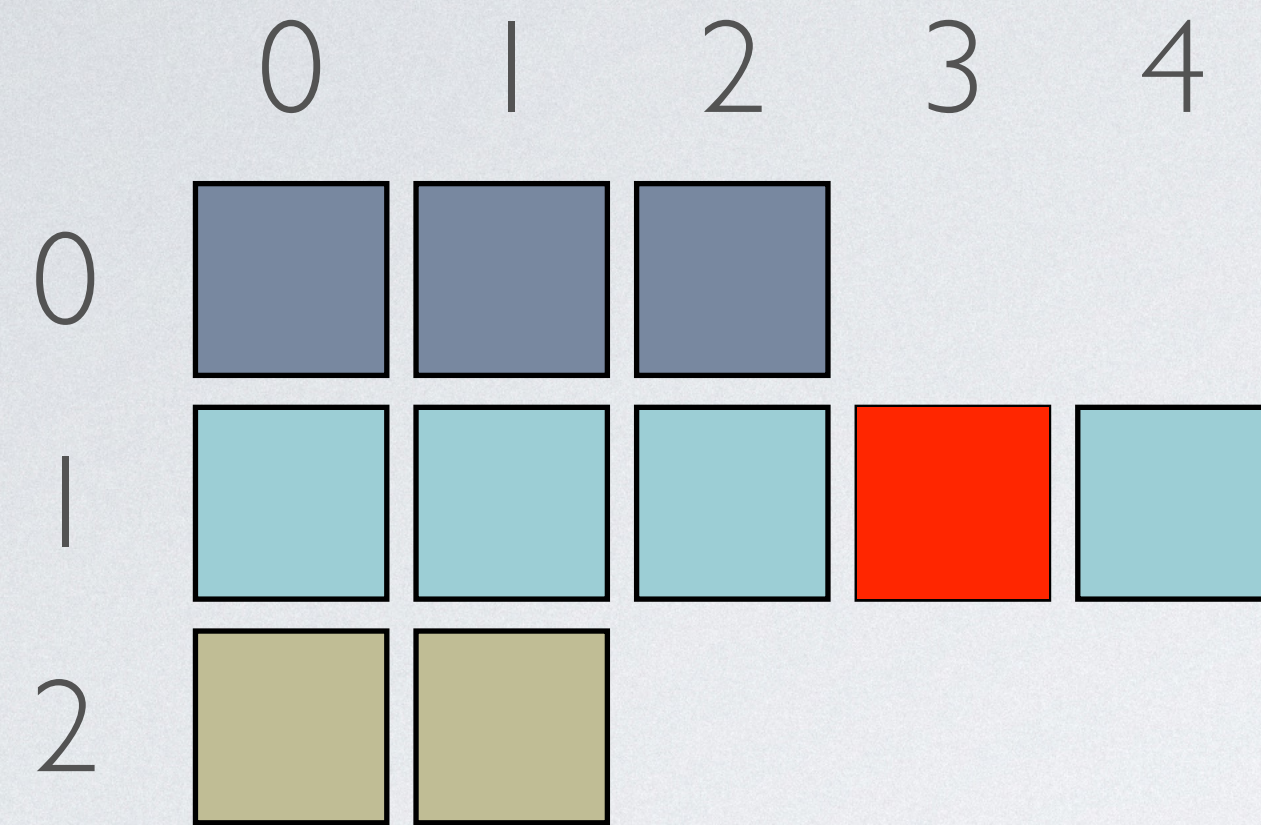
Generalize

- Compiler's loop representations
- Scheduling primitives and their impl.

Generalize

- Tensor storage scheme
- Tensor access lowering

Ragged Tensor Storage Without Padding



$$\text{Offset}(1, 3) = \text{RowStart}(1) + 3$$

Need to precompute dimension offsets before kernel execution

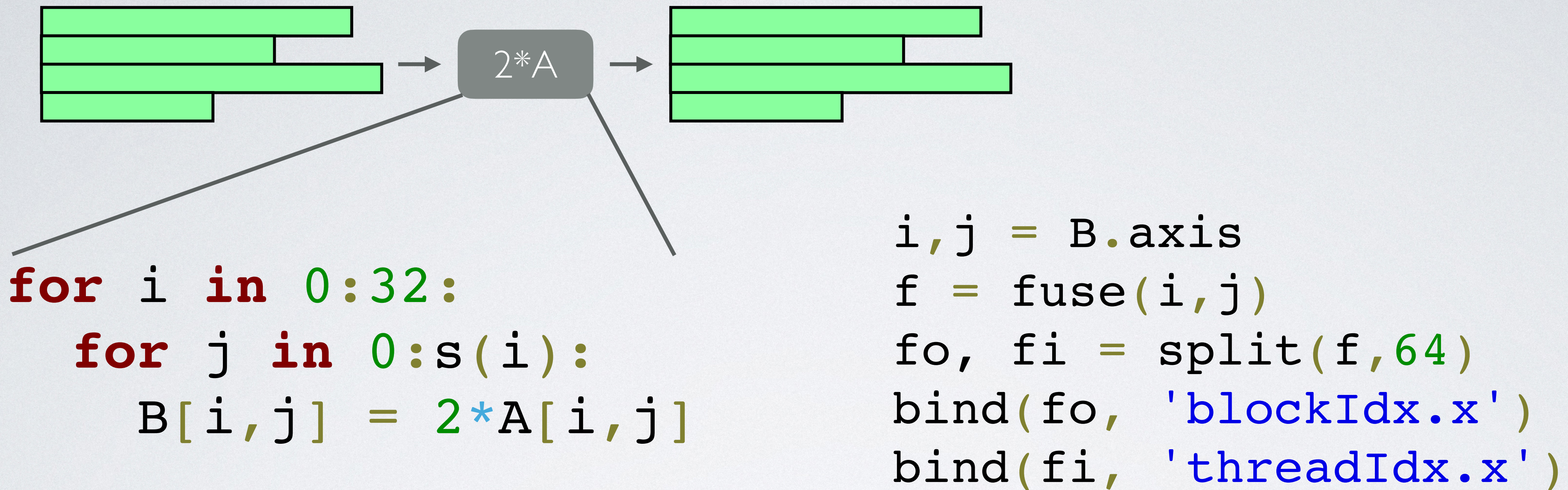


Once precomputed, we have cheap random accesses, similar to dense tensors!

Outline

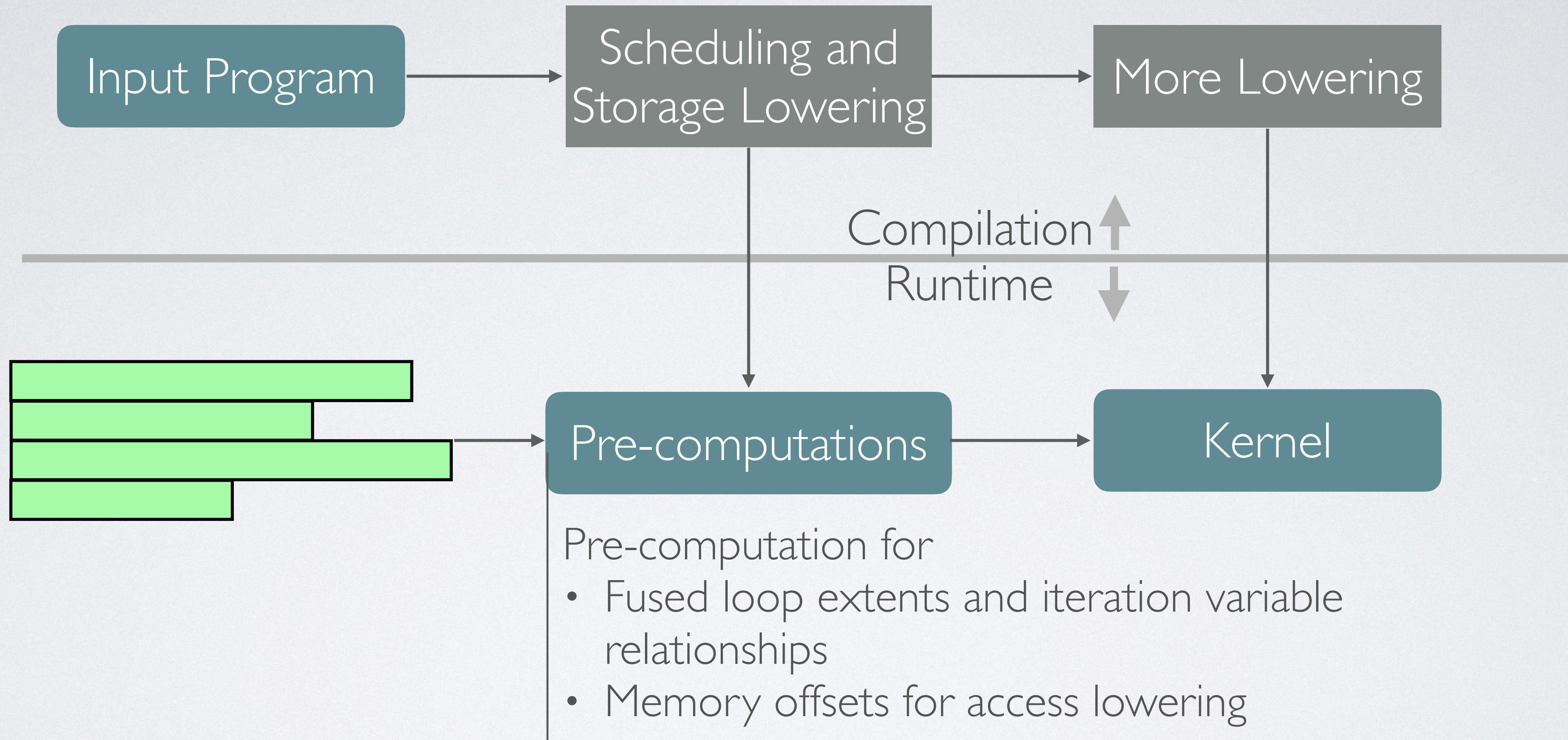
- Motivation: Inefficient Support for Ragged Tensors
- CoRa: Our Compiler Based Solution
 - Scheduling and lowering
 - **API and overview**
- Evaluation
- Wrapping up

CoRa's API Is Similar to That of Dense Compilers



Other scheduling primitives for load balancing, operation splitting, tensor dimension scheduling are available

CoRa's Compilation and Runtime Pipeline

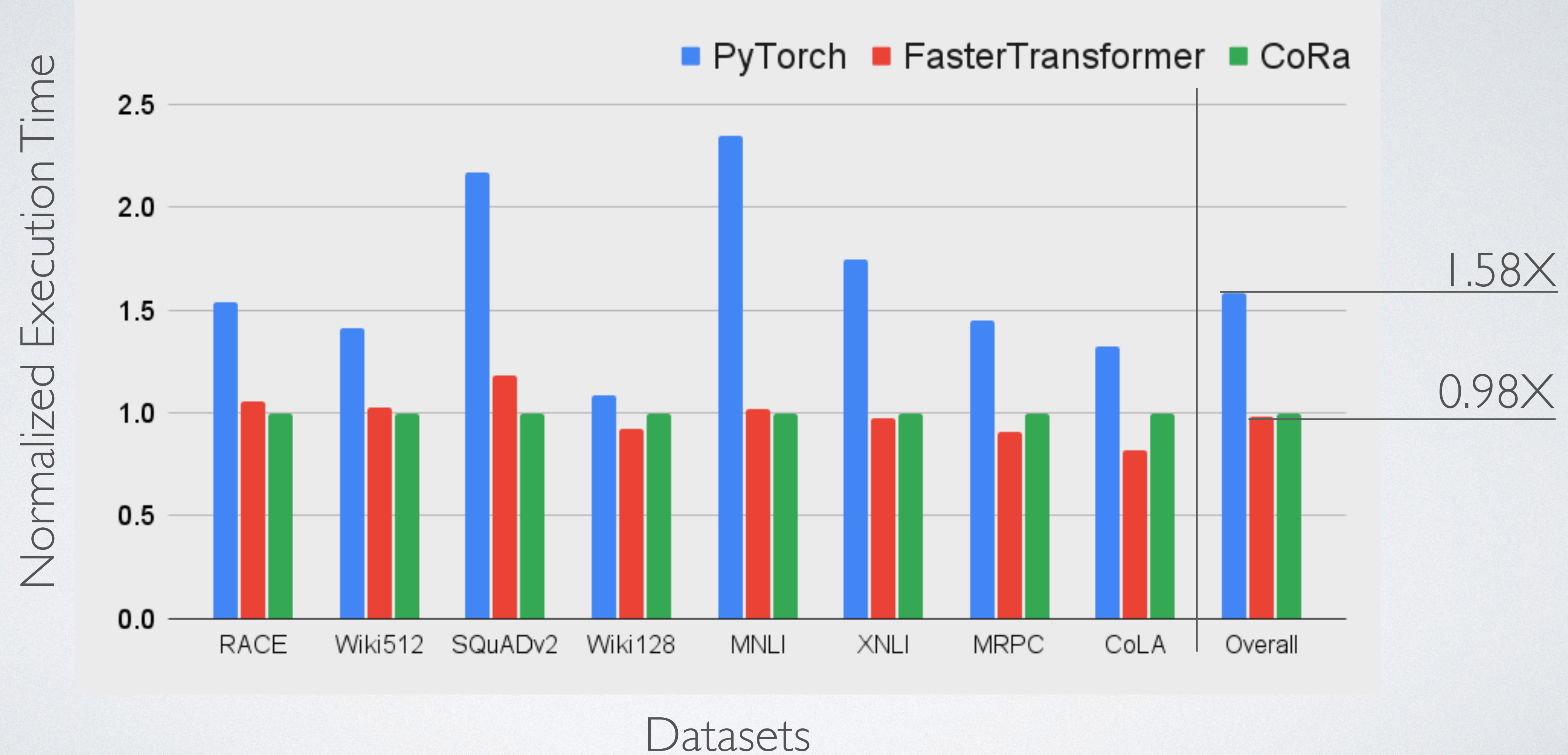


Outline

- Motivation: Inefficient Support for Ragged Tensors
- CoRa: Our Compiler Based Solution
 - Scheduling and lowering
 - API and overview
- **Evaluation**
- Wrapping up

Layer Forward Pass Latencies on Nvidia V100 GPU

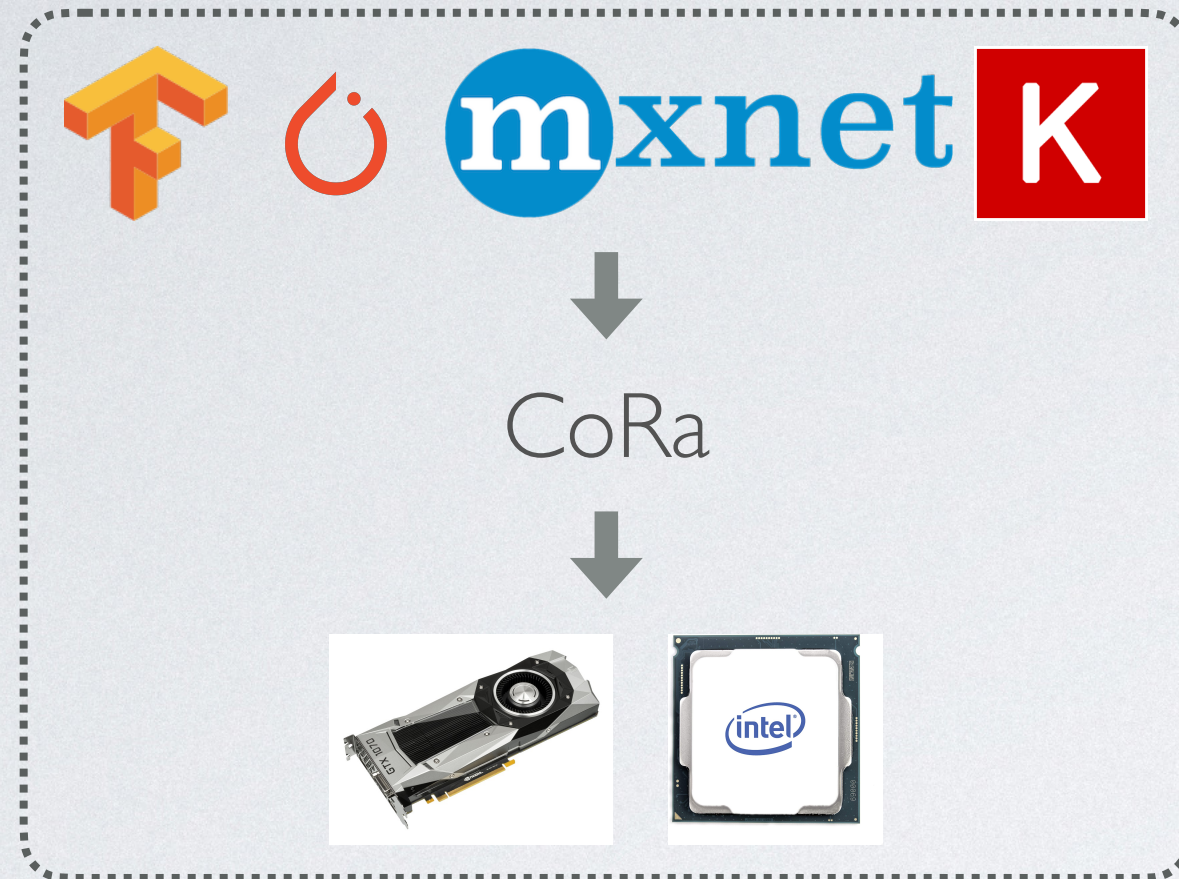
Lower is better



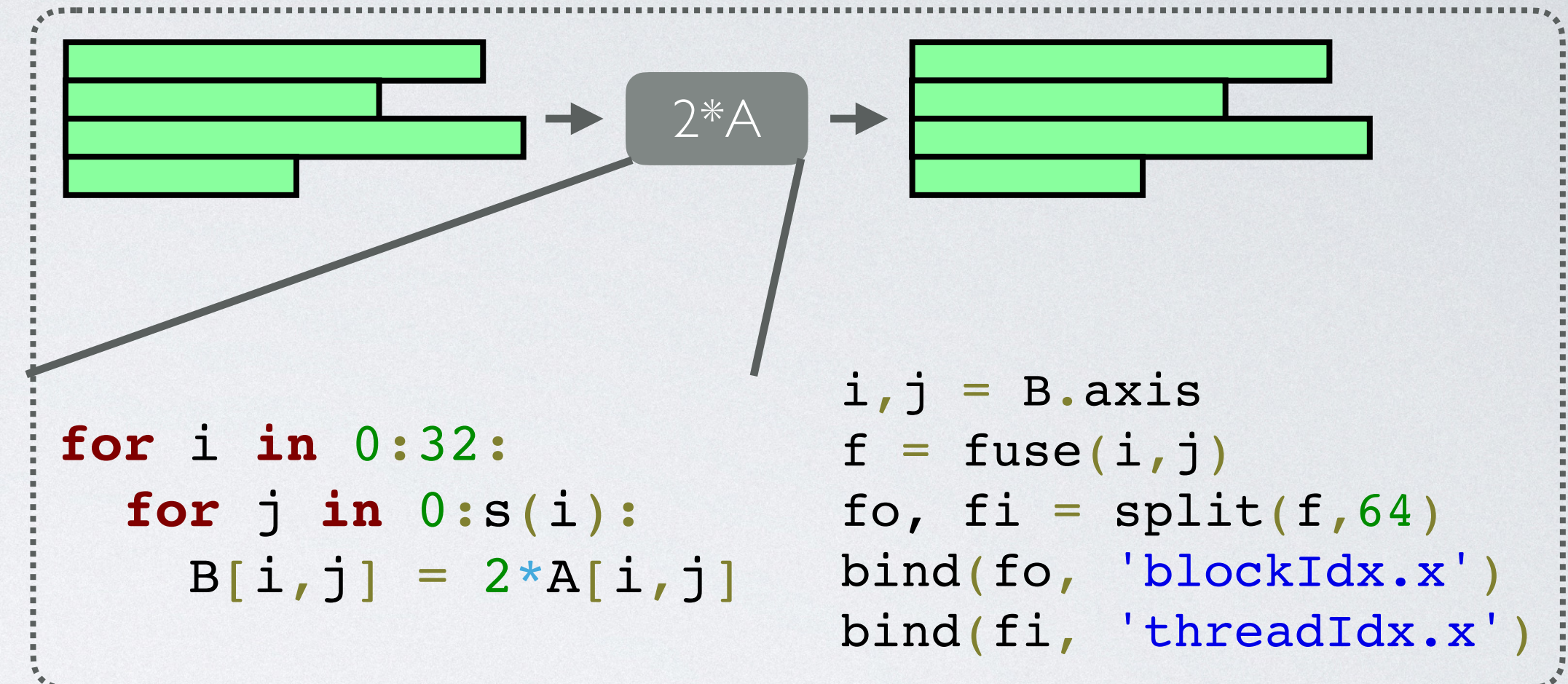
Outline

- Motivation: Inefficient Support for Ragged Tensors
- CoRa: Our Compiler Based Solution
 - Scheduling and lowering
 - API and overview
- Evaluation
- **Wrapping up**

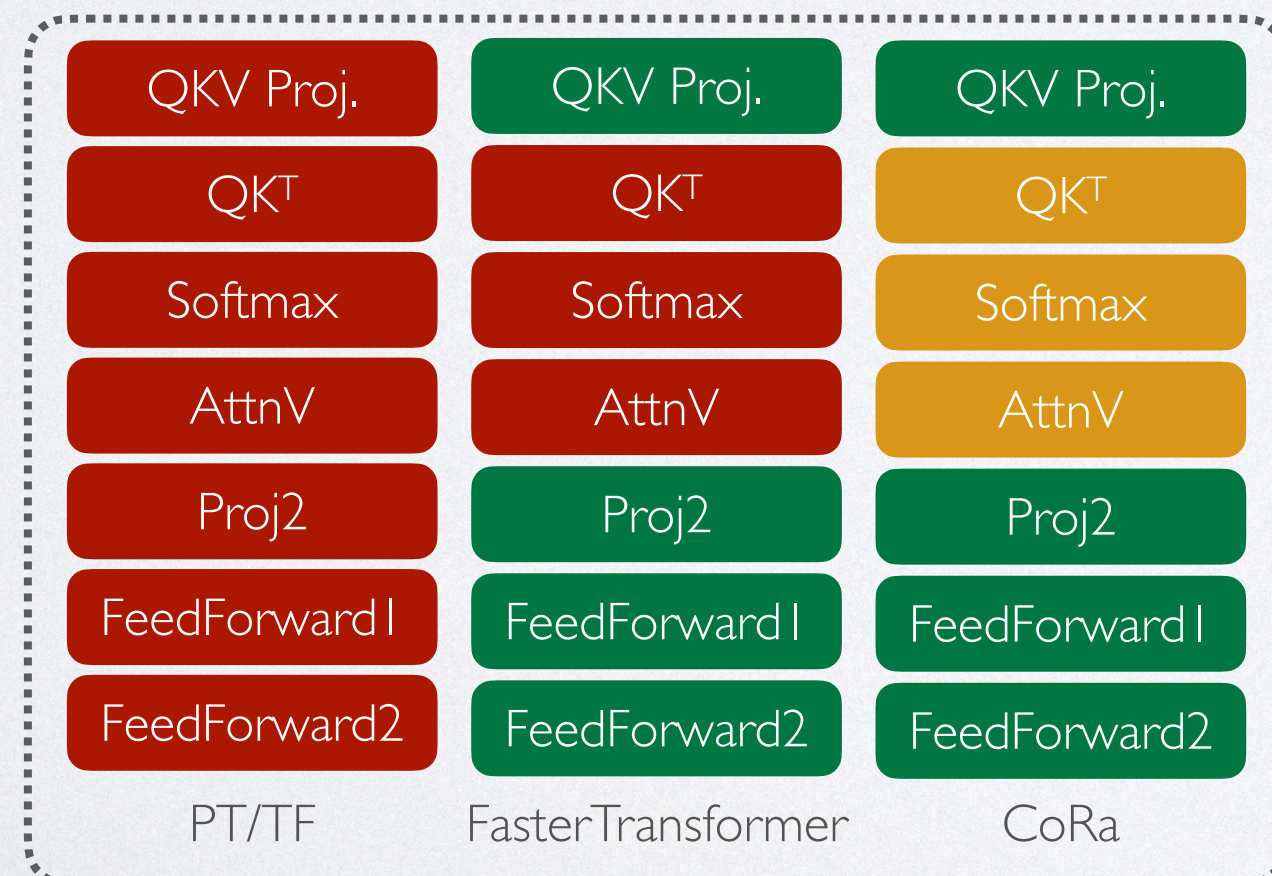
Wrapping Up CoRa



CoRa is a tensor compiler for operations on ragged tensors



CoRa provides a familiar API similar to that of dense tensor compilers



CoRa generates code as performant as hand-written code for transformer models