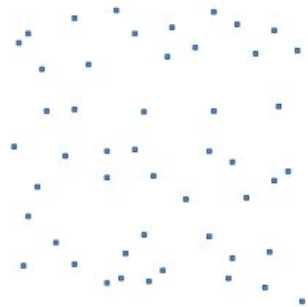
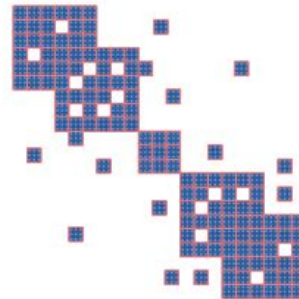
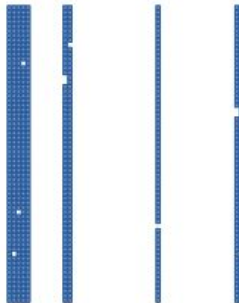


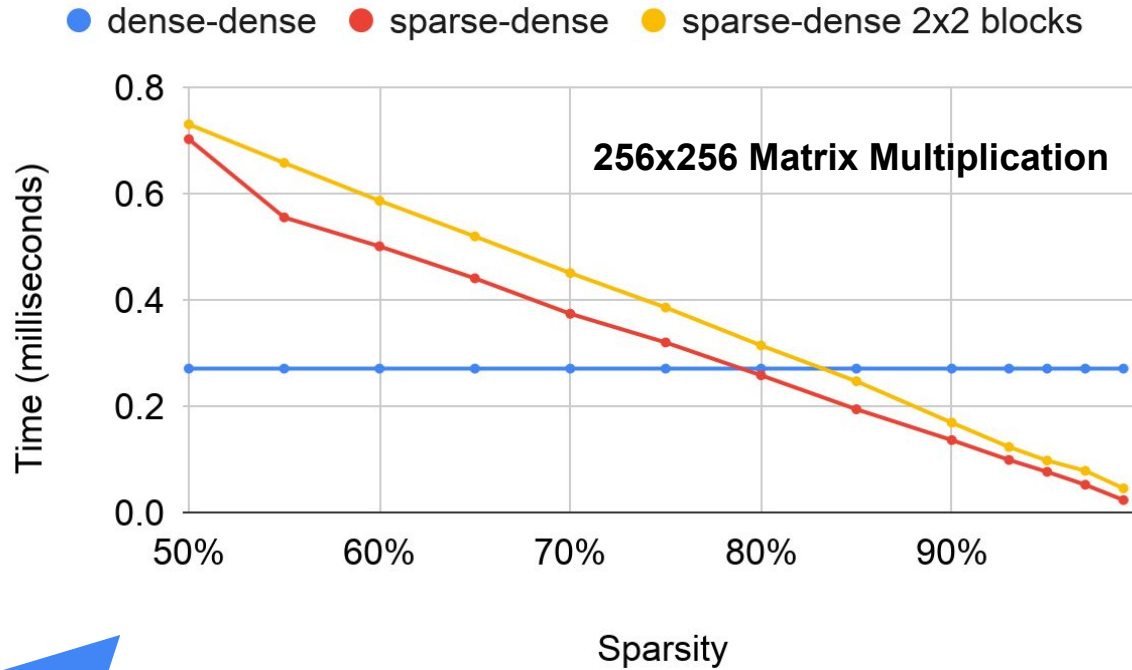
An ML-Driven Autoconfigurator for Sparse Tensor Kernels in MLIR

Gus Smith* w/ Aart Bik, Penporn Koanantakool, and
Mangpo Phothilimthana (Google)

*on internship from UW



sparsity takes many forms!



but, if we can exploit sparsity, it's immensely beneficial!

(from Penporn Koanantakool)

How do we exploit sparsity?

With sparsity formats!

Indices:	(0, 0)	(0, 2)	(0, 3)	(2, 0)	(2, 3)
Values:	6	9	8	5	7

Coordinate (COO)

Row pointers:	0	3	3	5	
Column indices:	0	2	3	0	3
Values:	6	9	8	5	7

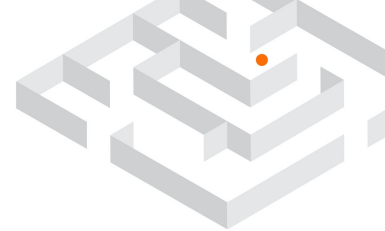
Compressed Sparse Row (CSR)

Column pointers:	0	2	2	3	5
Row indices:	0	2	0	0	2
Values:	6	5	9	8	7

Compressed Sparse Column (CSC)

plus DIA, ELL, and more...

...and then along came TACO...

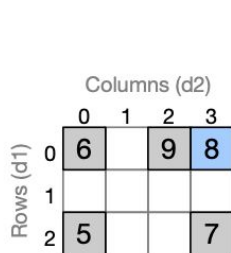


Dimension-wise Specification

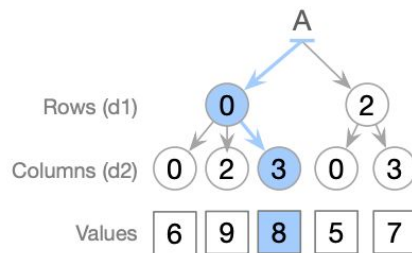
- Each dimension can be sparse or dense
- Define traversal order: (1st dim, 2nd dim, ...)
- For a k^{th} -order tensor, this covers $k!2^k$ formats

Dense: All elements in that dimension are stored (e.g., dense row pointers in CSR).

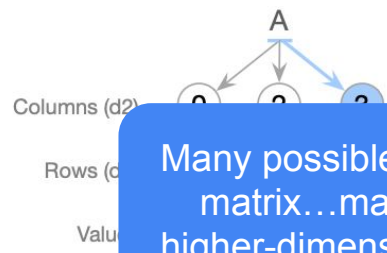
Sparse: Zeros in that dimension aren't stored.



(a) Sparse matrix A



(b) Row-major coordinate storage tree



(c) Column-major coordinate storage tree

Many possible formats for a matrix...many more for higher-dimensional tensors!

size $\boxed{3}$
Dense
Rowmajor

size $\boxed{4}$
vals $\boxed{6\ 0\ 9\ 8\ 0\ 0\ 0\ 0\ 5\ 0\ 0\ 7}$

(d) dense_{d1},dense_{d2}

pos $\boxed{0\ 2}$
idx $\boxed{0\ 2}$
size $\boxed{4}$ Row-slicing
vals $\boxed{6\ 0\ 9\ 8\ 5\ 0\ 0\ 7}$

(e) sparse_{d1},dense_{d2}

size $\boxed{3}$ CSR
pos $\boxed{0\ 3\ 3\ 5}$
idx $\boxed{0\ 2\ 3\ 0\ 3}$
vals $\boxed{6\ 9\ 8\ 5\ 7}$

(f) dense_{d1},sparse_{d2}

pos $\boxed{0\ 2}$
idx $\boxed{0\ 2}$ DCSR
pos $\boxed{0\ 3\ 5}$
idx $\boxed{0\ 2\ 3\ 0\ 3}$
vals $\boxed{6\ 9\ 8\ 5\ 7}$

(g) sparse_{d1},sparse_{d2}

size $\boxed{4}$
Dense
Colmajor

size $\boxed{3}$
vals $\boxed{6\ 0\ 5\ 0\ 0\ 0\ 9\ 0\ 0\ 8\ 0\ 7}$

(h) dense_{d2},dense_{d1}

pos $\boxed{0\ 3}$
idx $\boxed{0\ 2\ 3}$
size $\boxed{3}$ Col-slicing
vals $\boxed{6\ 0\ 5\ 9\ 0\ 0\ 8\ 0\ 7}$

(i) sparse_{d2},dense_{d1}

size $\boxed{4}$ CSC
pos $\boxed{0\ 2\ 2\ 3\ 5}$
idx $\boxed{0\ 2\ 0\ 0\ 2}$
vals $\boxed{6\ 5\ 9\ 8\ 7}$

(j) dense_{d2},sparse_{d1}

pos $\boxed{0\ 3}$
idx $\boxed{0\ 2\ 3}$ DCSC
pos $\boxed{0\ 2\ 3\ 5}$
idx $\boxed{0\ 2\ 0\ 0\ 2}$
vals $\boxed{6\ 5\ 9\ 8\ 7}$

(k) sparse_{d2},sparse_{d1}

thanks to Aart and others, we can now easily configure sparse kernels in MLIR!

```
#CSR = #sparse_tensor.encoding<{  
  dimLevelType = [ "dense", "compressed" ],  
  dimOrdering = affine_map<(d0, d1) -> (d0, d1)>  

```

```
func @kernel(%a: tensor<?x?xf64, #CSR>,  
             %b: tensor<?x?xf64>,  
             %c: tensor<?x?xf64>) -> tensor<?x?xf64> {  
  
  ...kernel implementation...  
  
  return %d : tensor<?x?xf64>  
}
```

Goal: given a kernel, can we
configure its sparse format?



couldn't you just brute force it?

Goal: given **any** kernel, can we
quickly configure its sparse
format?

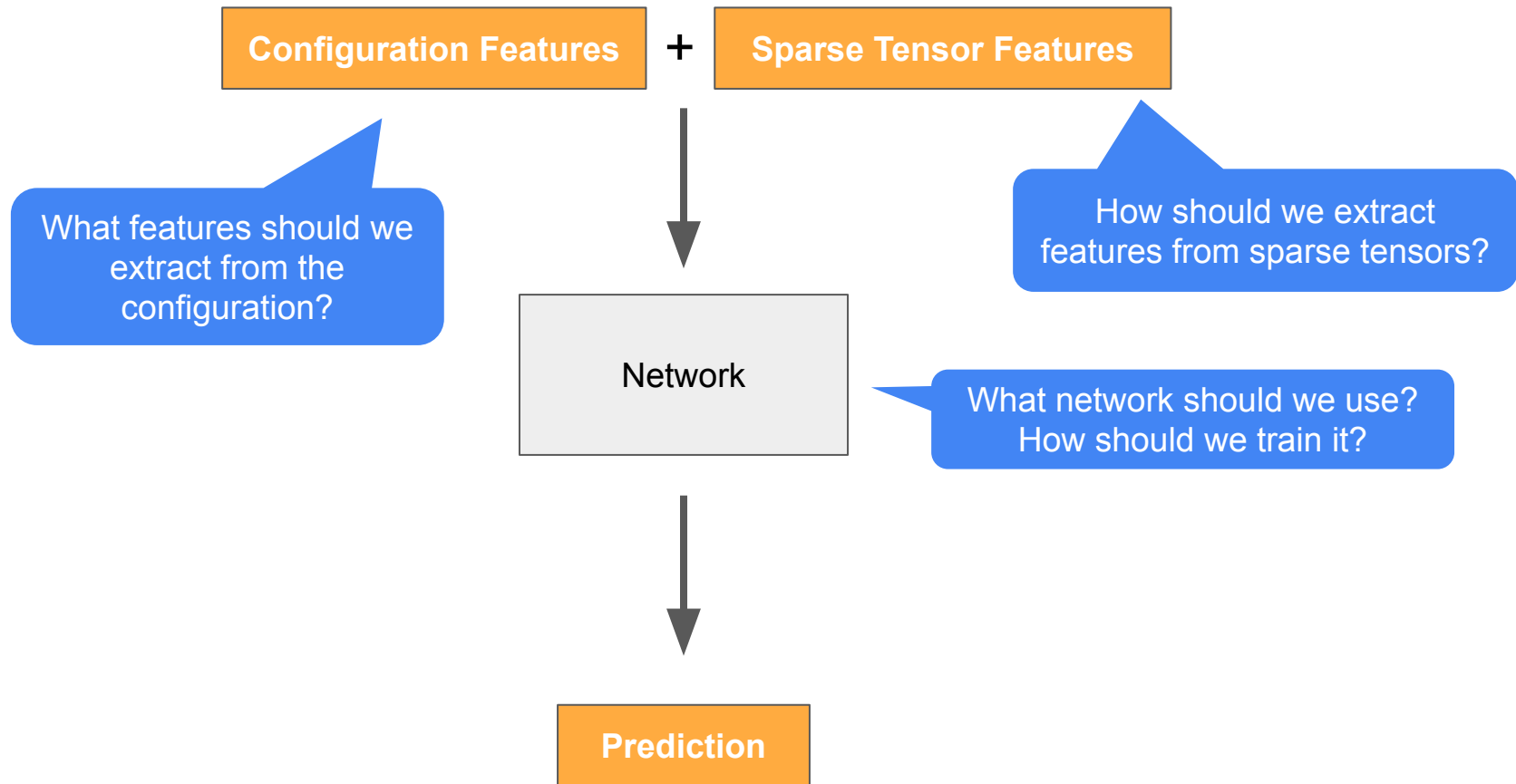
so we could use it in a JIT!

We use a “standard” ML-for-ML approach:

1. Train a cost model
2. Use the cost model in a search procedure

(A Learned Performance Model for Tensor Processing Units (Kaufman et. al.)
Learning to Optimize Halide with Tree Search and Random Programs (Adams et. al.)
Learning to Optimize Tensor Programs (Chen et. al.))

Training a Cost Model



Configuration Features

What features should we extract from the configuration?

Configuration Features

- Entire sparse kernel configuration is easily featurized by packing into fixed-length vector
- Per-dimension sparse formats, dimension ordering, plus other settings e.g. parallelization and vectorization levels

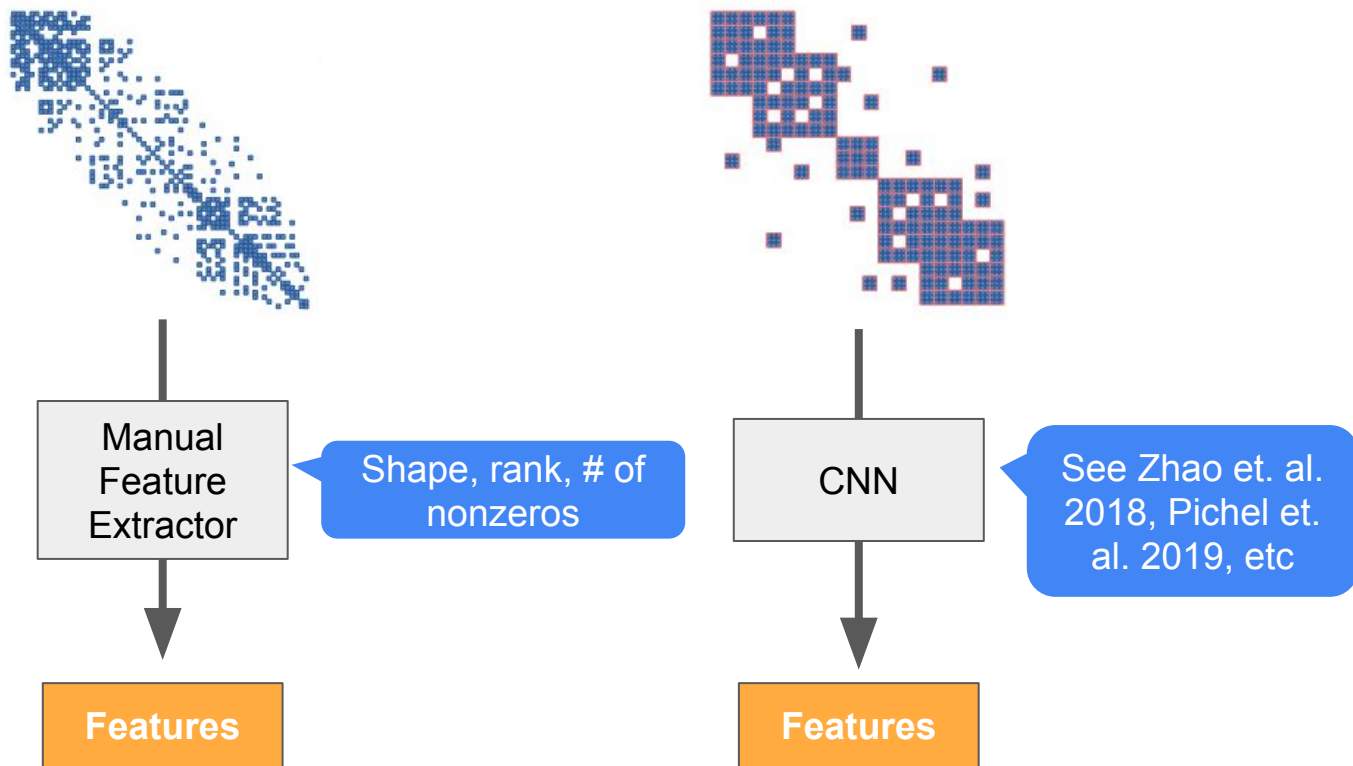
For sparse-dense matmul
kernel: ~9k configurations

Configuration Features

+

Sparse Tensor Features

How should we extract
features from sparse tensors?



Bridging the Gap between Deep Learning and Sparse Matrix Format Selection, Zhao et. al 2018

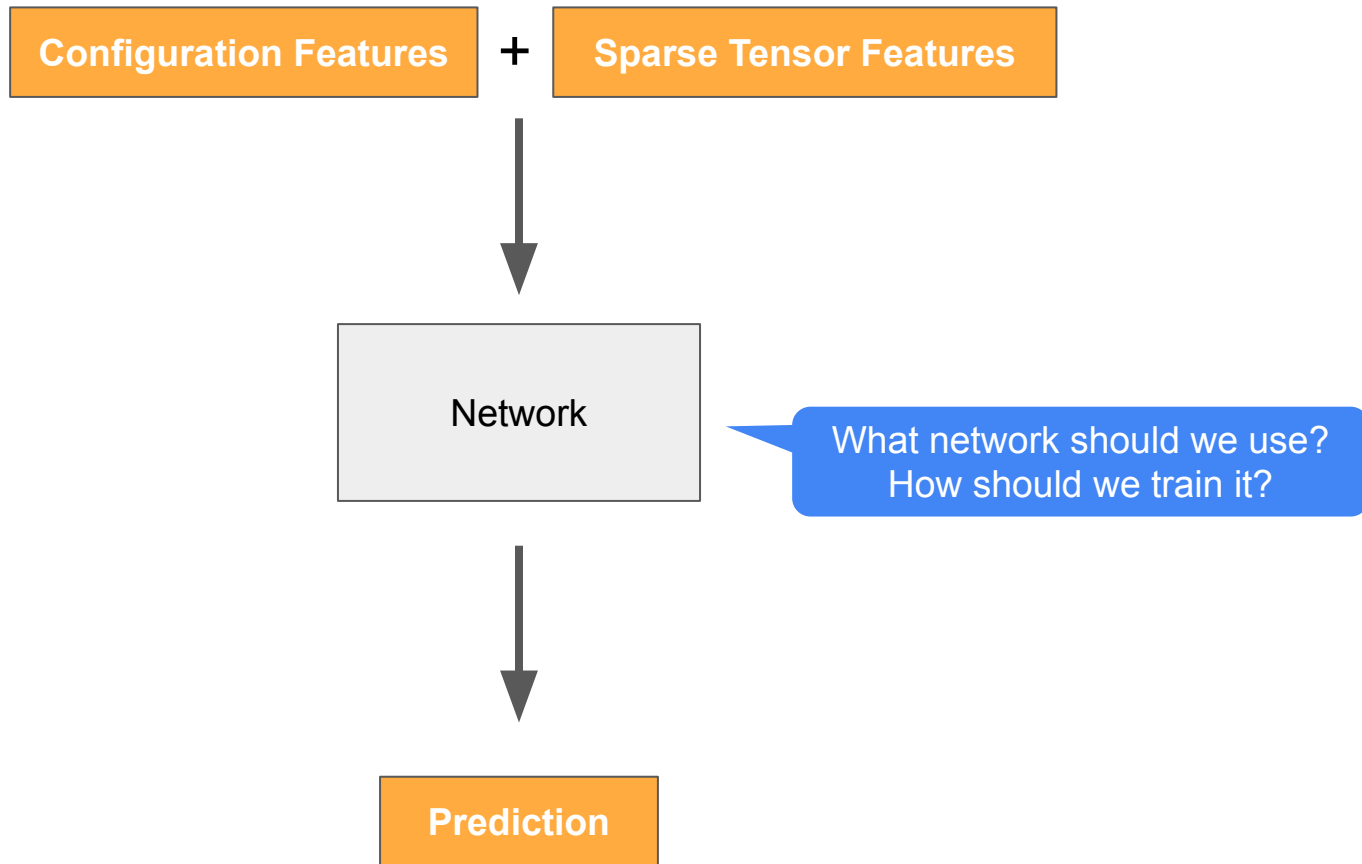
Sparse Matrix Classification on Imbalanced Datasets Using Convolutional Neural Networks, Pichel et. al. 2019

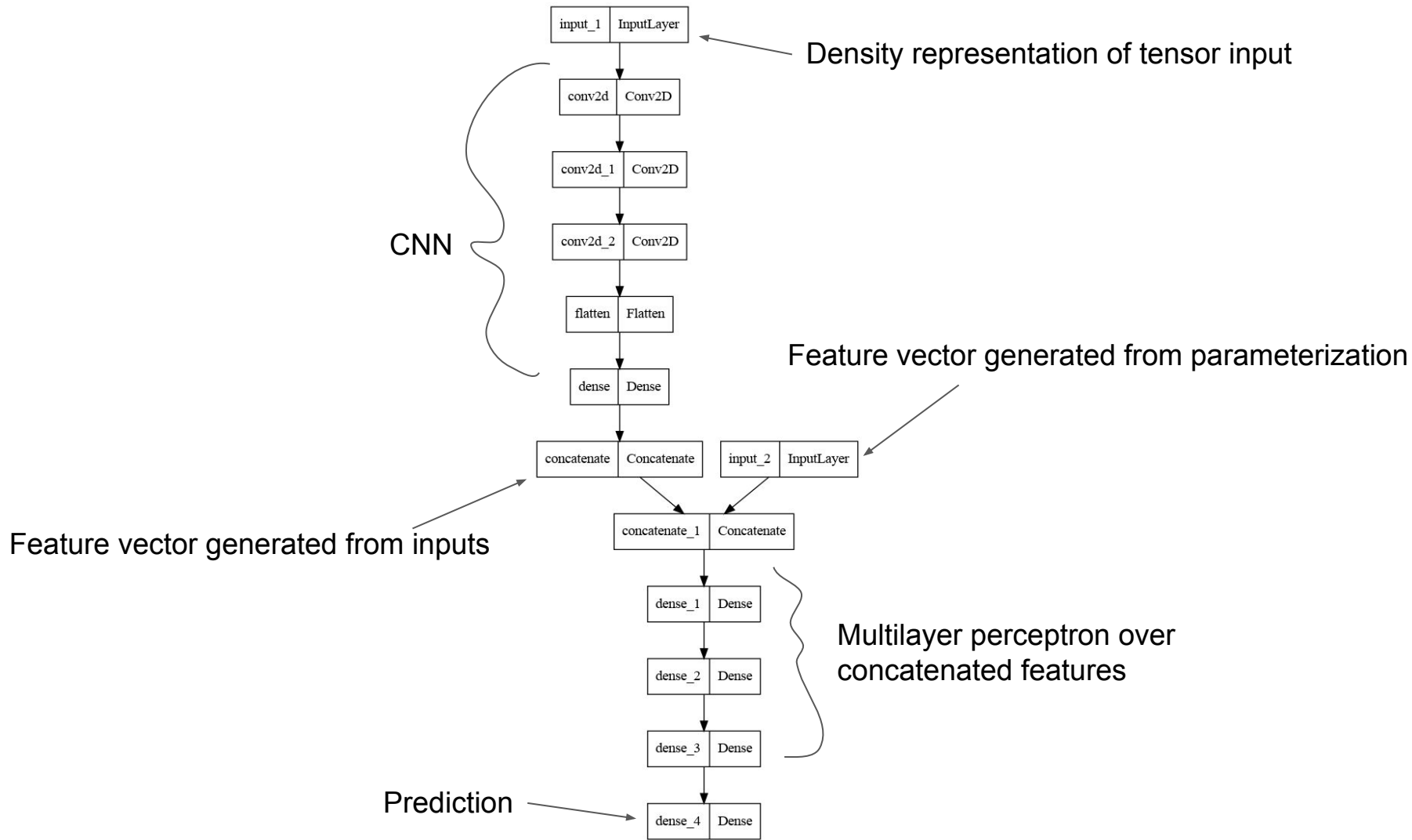
IA-SpGEMM: an input-aware auto-tuning framework for parallel sparse matrix-matrix multiplication, Xie et. al. 2019

Variable-sized **sparse** input matrix

1	2			3	4		
5	6						
		7	8				
		9	10				
				11	12		
13	14			15	16		
						17	18
						19	20

Now we can perform dense 2D convolution!





9k parameterizations

X

3k tensor inputs

= 27M points!

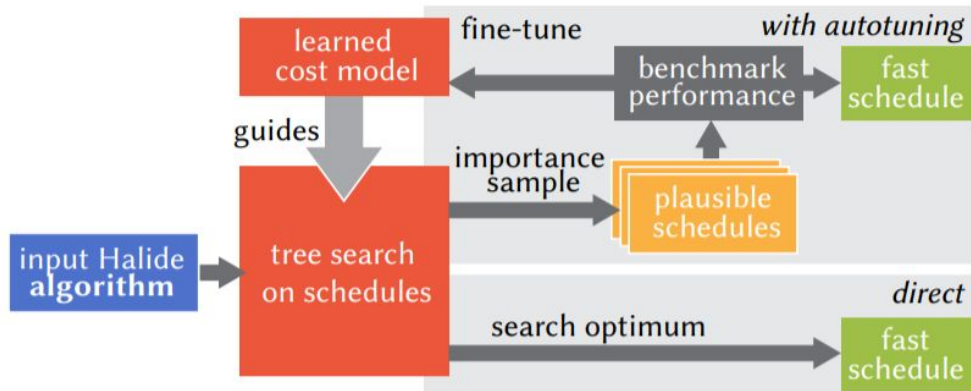
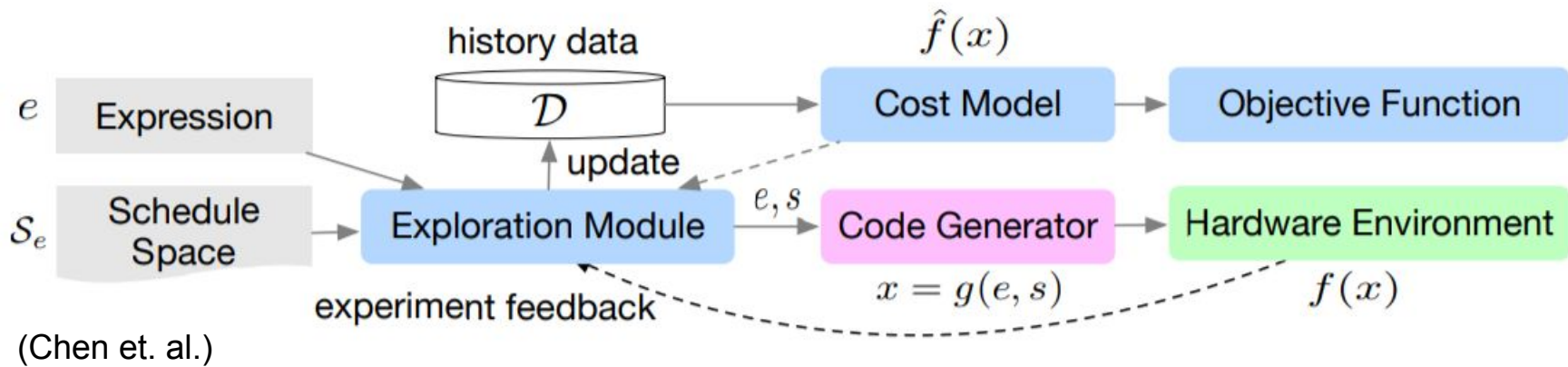
From sparse ResNet50 dataset
(Sparse GPU Kernels for Deep
Learning, Gale et. al. 2020)

Randomly sample and benchmark the runtimes of
50k (parameterization, input) points

Benchmarked on an Intel Xeon
6154 (Skylake), 72 CPUs

```
model.compile(  
    optimizer='adam',  
    loss='mean_squared_error')  
model.fit(...)
```

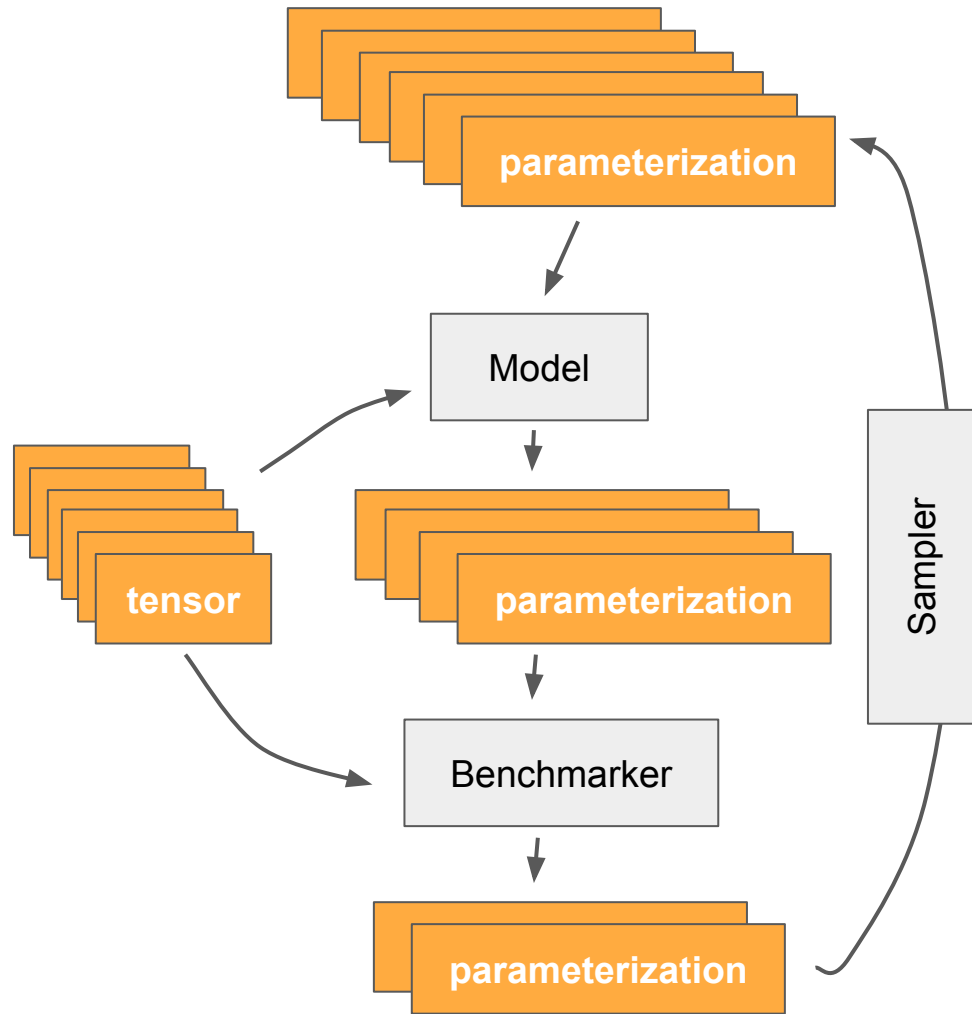
Search

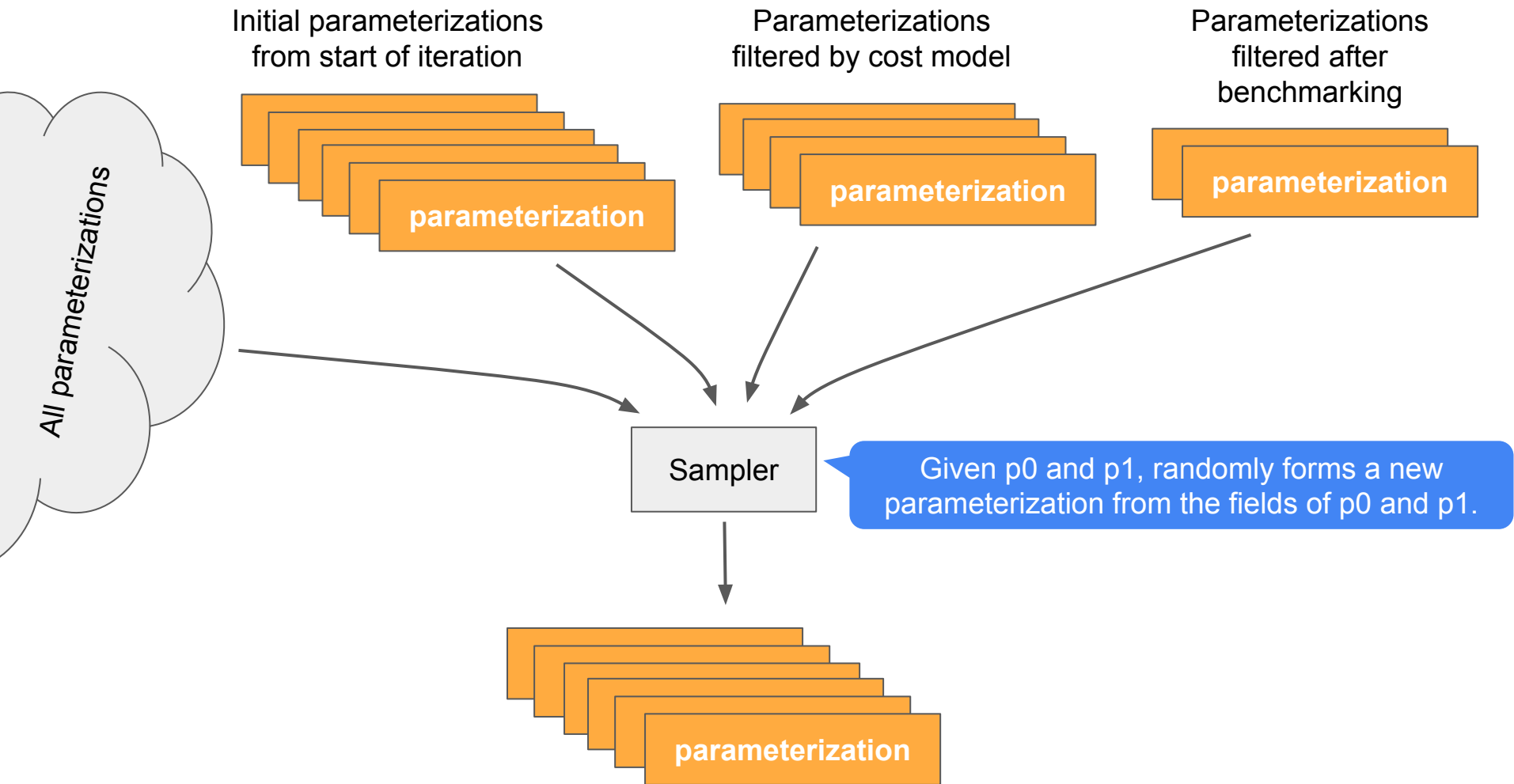


(Adams et. al.)

In related works, search follows similar structure

1. Randomly select test inputs
2. Randomly sample initial parameterizations
3. For each search round:
 - a. Use the cost model to filter the set of candidates (fast)
 - b. Benchmark the filtered set of candidates (slow)
 - c. Record top benchmarked candidates
 - d. Construct next generation of parameterizations





Future work

- Parameterize cost model over kernels
 - Generate features from kernels, e.g. via manual feature extraction or GNN
- Train on data from multiple kernels (SDDMM; random sparse kernels)

Thank you!