

SSketch: An Automated Framework for Streaming Sketch-based Analysis of Big Data on FPGA*

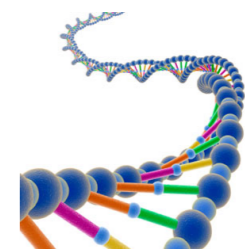
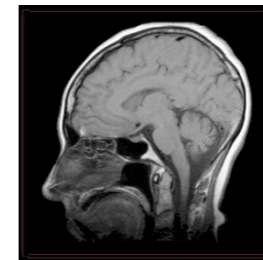
Bitva Darvish Rouhani, Ebrahim Songhori, Azalia Mirhoseini, and Farinaz Koushanfar
Department of ECE, Rice University
Houston Texas, USA



* This work was supported in parts by the Office of Naval Research grant (ONR- N00014-11-1-0885)

Introduction

- Era of big data
 - Decision making
 - Find patterns
 - Prevent failures
- Applications
 - Medical
 - Cybersecurity
 - Media/social networks
 - Finance
- Machine learning and statistical optimizations are main enablers



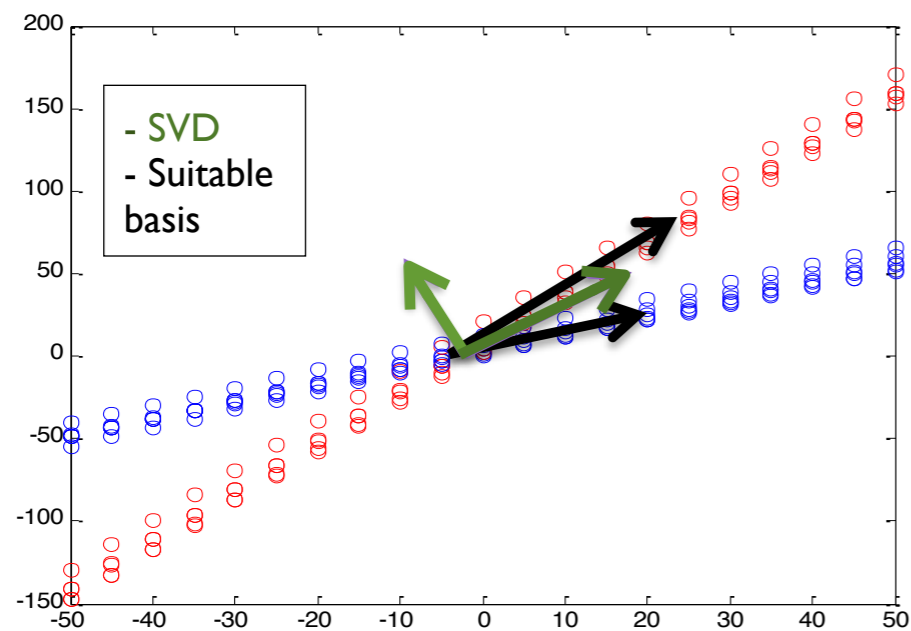
Efficient Data Transformation

- Data transformation
 - Compact representation of the data collection
 - Exploiting the redundancy present in the dataset
- An efficient data transformation should simultaneously consider the:
 - Scalability
 - Application
 - Underlying platform constraints



Ensemble of Lower Dimensional Structures

- Many modern massive datasets are either low-rank or lie on a union of lower dimensional subspaces^[1]



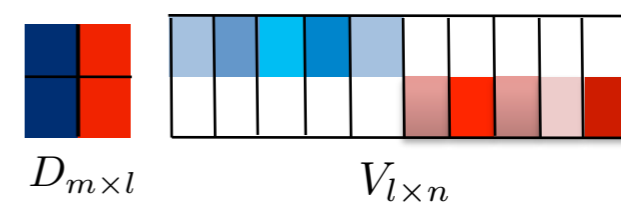
Original dense dataset $A_{m \times n}$



Singular value decomposition on A



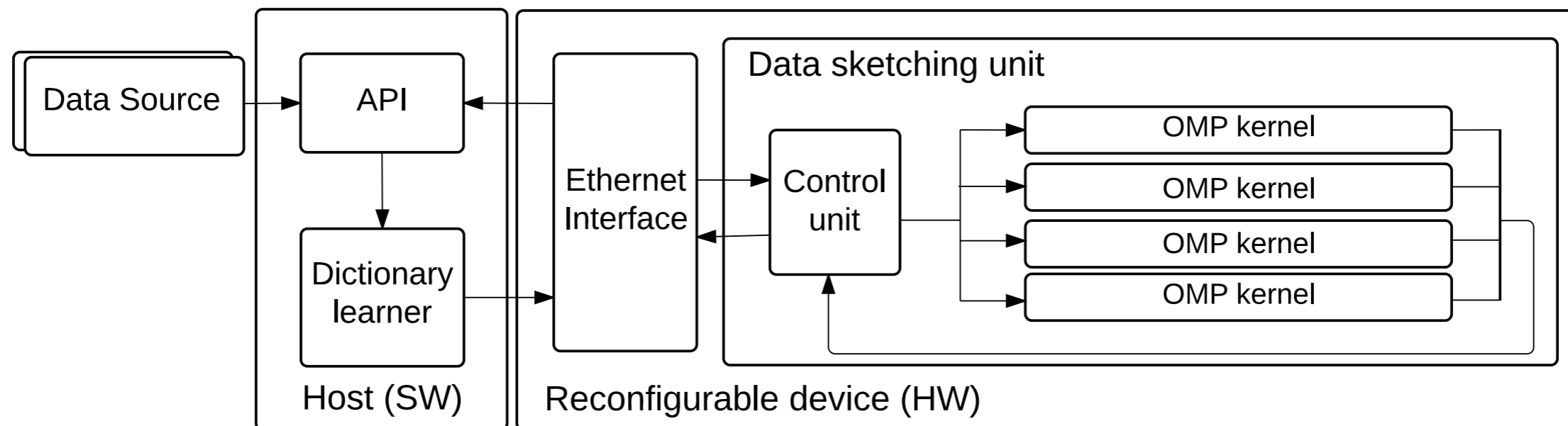
Ensemble of lower dimensional structures



[1] Mirhoseini et. al, "Rankmap: A Platform-aware Framework for Distributed Learning from Dense Datasets", arXiv:1503.08169, 2015

SSketch Framework

- Streaming Sketch-based analysis of big data using FPGA
- Transforming the big data with dense correlations to an ensemble of lower dimensional subspaces
- Streaming applications:
 - Limited storage
 - Single pass access to data



SSketch Framework

- The transformed data is applicable to a broad set of matrix-based data analysis algorithms:
 - Regularized loss function optimization
 - Power method
 - Image processing
 - De-noising
 - Super-resolution
 - Classification



SSketch Methodology

- Data transformation:

$$\underset{\mathbf{D} \in R^{m \times l}, \mathbf{V} \in R^{l \times n}}{\text{minimize}} \quad \|\mathbf{A} - \mathbf{D}\mathbf{V}\|_F \quad \text{subject to} \quad \|\mathbf{V}\|_0 \leq kn$$

- Adaptive error-based dictionary learning
- Single pass access to data

Algorithm 1 SSketch algorithm

Inputs: Measurement matrix \mathbf{A} , projection threshold α , sparsity level k , error threshold ϵ , and dictionary size l .

Output: Matrix \mathbf{D} , and coefficient matrix \mathbf{V} .

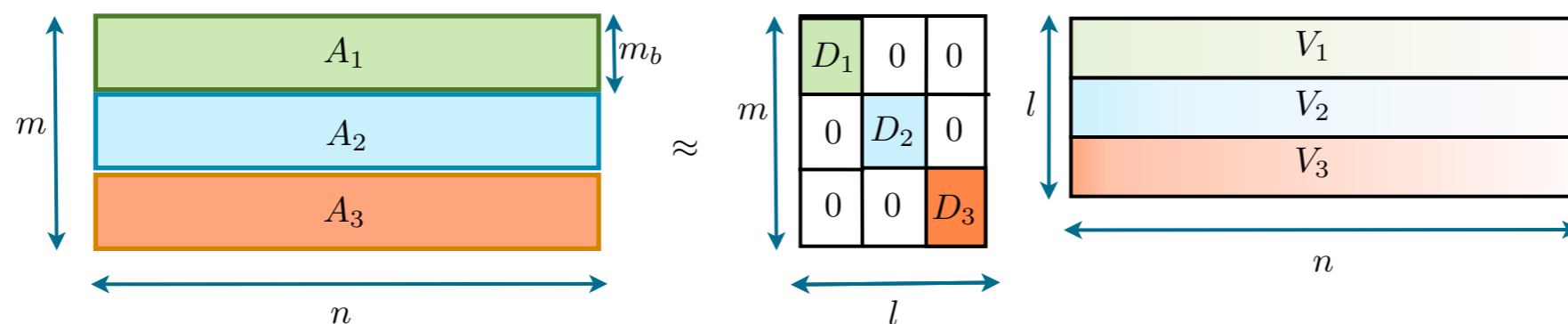
```
1:  $\mathbf{D} \leftarrow$  empty
2:  $j \leftarrow 0$ 
3: for  $i = 1, \dots, n$  do
4:    $W(\mathbf{A}_i) = \frac{\|\mathbf{D}(\mathbf{D}^t\mathbf{D})^{-1}\mathbf{D}^t\mathbf{A}_i - \mathbf{A}_i\|_2}{\|\mathbf{A}_i\|_2}$ 
5:   if  $W(\mathbf{A}_i) > \alpha$  and  $j < l$  then
6:      $\mathbf{D}_j = \mathbf{A}_i / \sqrt{\|\mathbf{A}_i\|_2}$ 
7:      $\mathbf{V}_{ij} = \sqrt{\|\mathbf{A}_i\|_2}$ 
8:      $j \leftarrow j + 1$ 
9:   else
10:     $\mathbf{V}_i \leftarrow \text{OMP}(\mathbf{D}, \mathbf{A}_i, k, \epsilon)$ 
11:  end if
12: end for
```

SSketch Methodology

- Data transformation:

$$\underset{\mathbf{D} \in R^{m \times l}, \mathbf{V} \in R^{l \times n}}{\text{minimize}} \quad \|\mathbf{A} - \mathbf{D}\mathbf{V}\|_F \quad \text{subject to} \quad \|\mathbf{V}\|_0 \leq kn$$

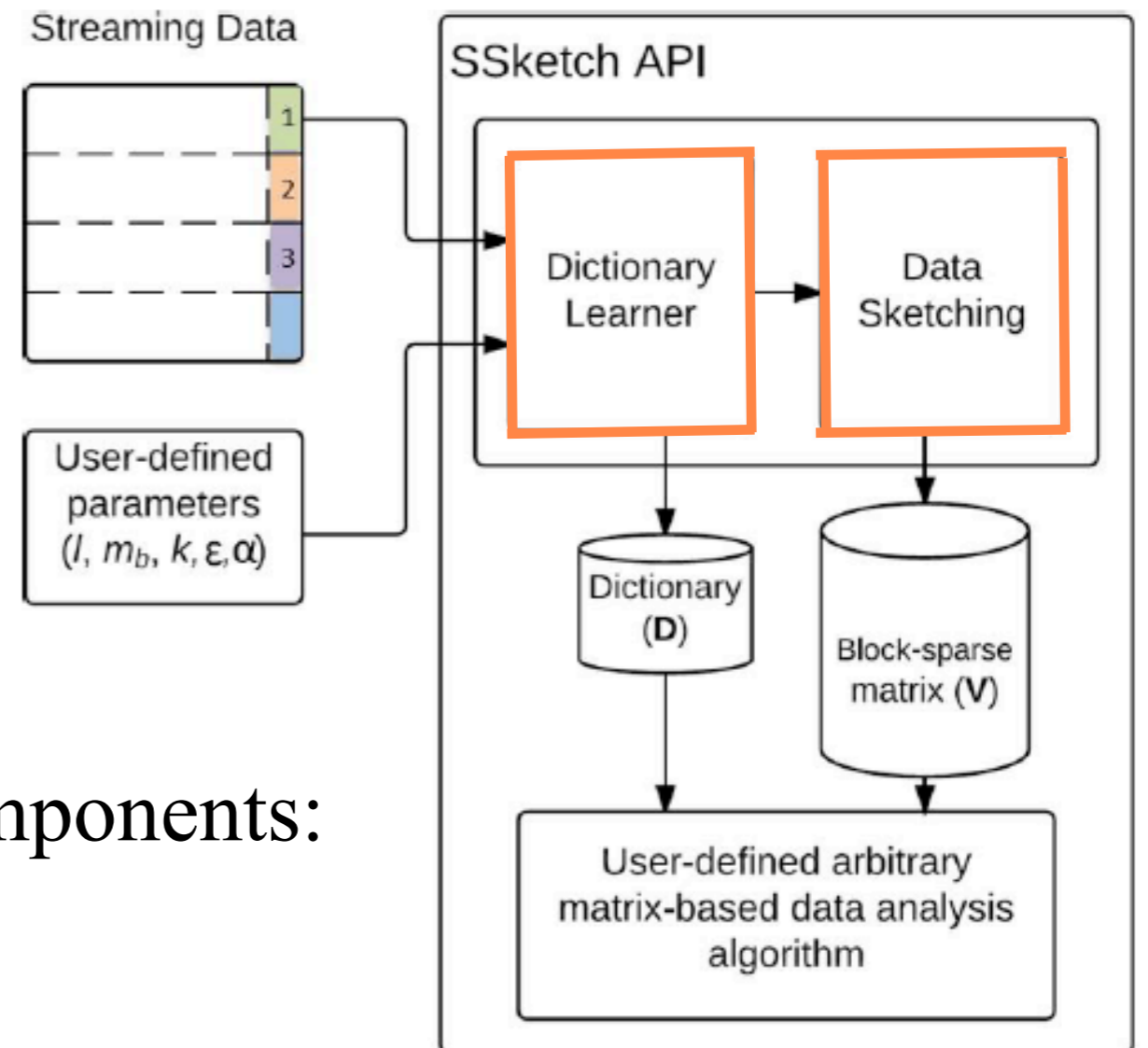
- Block splitting
- Amenable to FPGA accelerator



Schematic depiction of blocking SSketch

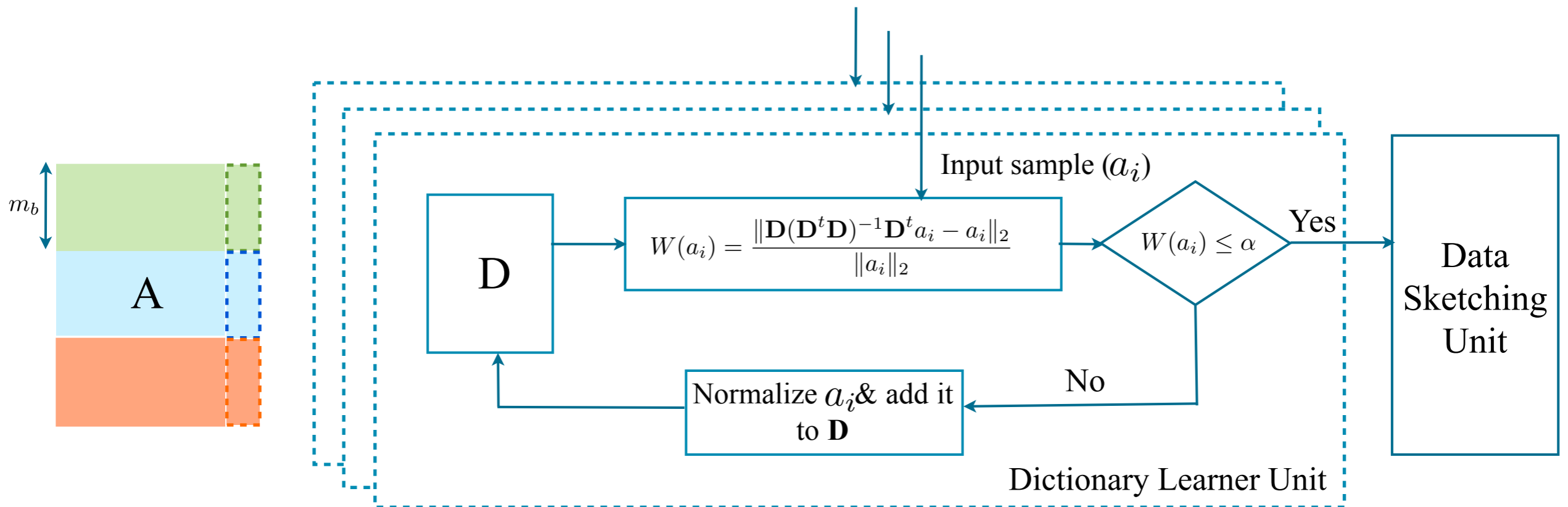
SSketch API

- Inputs:
 - Stream of input data
 - SSketch Algorithmic parameters
- Outputs:
 - Dictionary matrix D
 - Block-sparse matrix V
- SSketch consists of two main components:
 - Dictionary learner unit
 - Data sketching unit



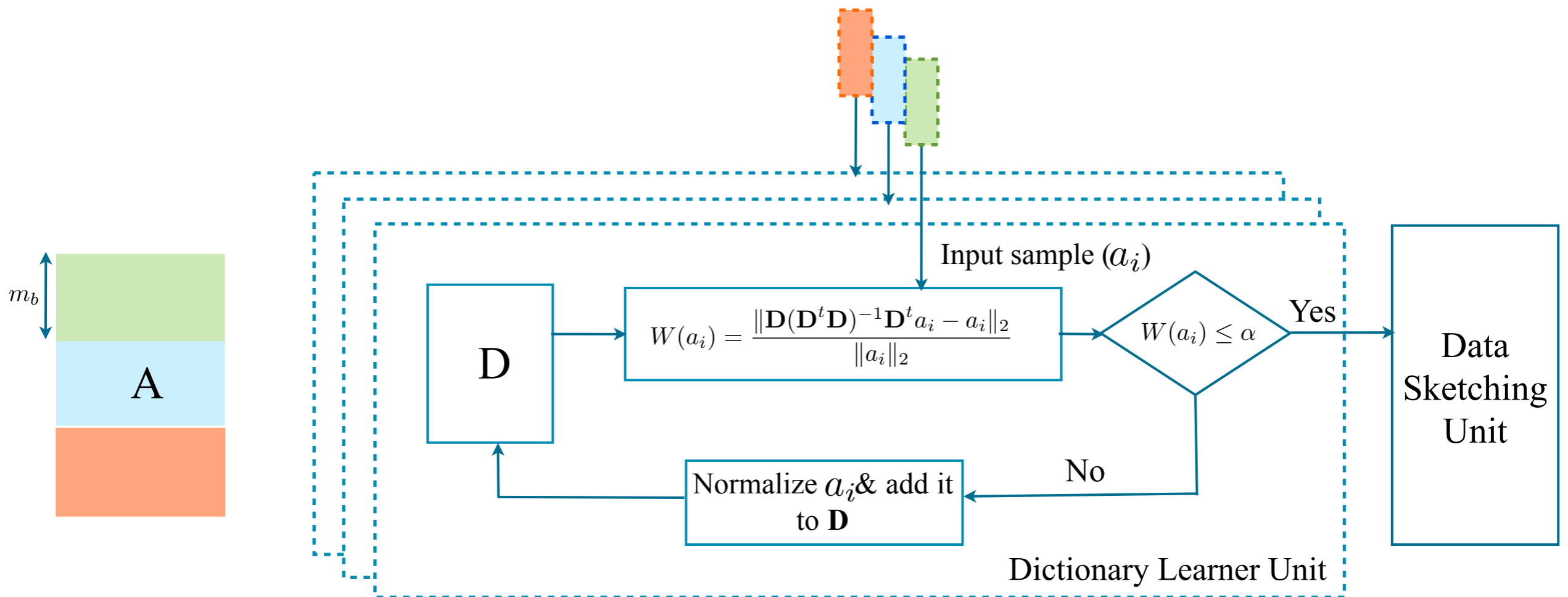
Adaptive Dictionary Learning

- By Streaming the input data, SSketch adaptively:
 - Learns/updates the corresponding dictionary of each block
 - Computes the data sketch



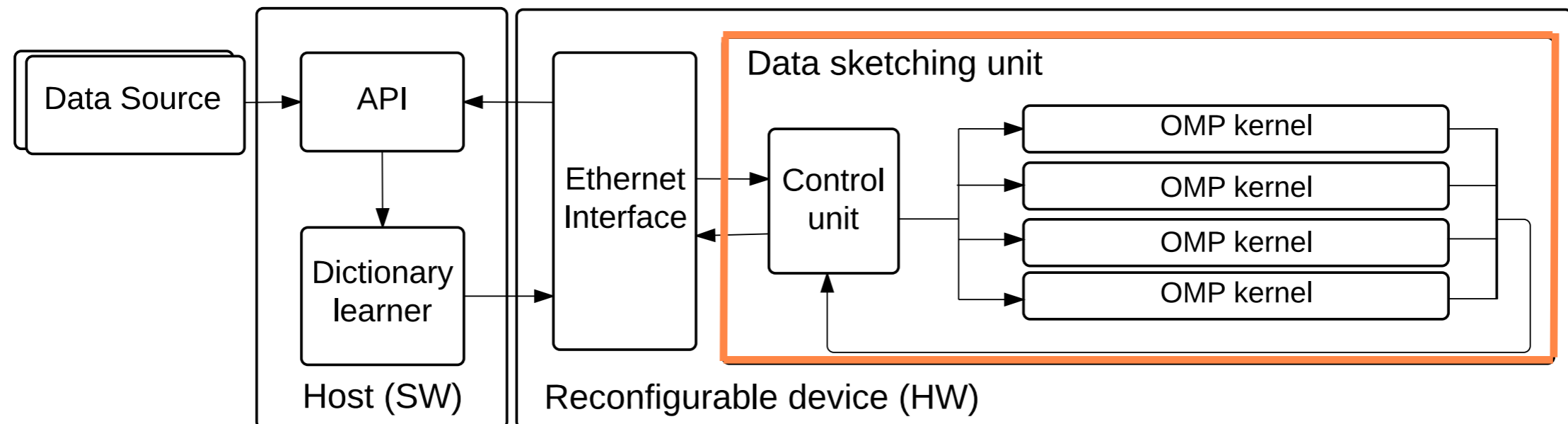
Adaptive Dictionary Learning

- By Streaming the input data, SSketch adaptively:
 - Learns/updates the corresponding dictionary of each block
 - Computes the data sketch



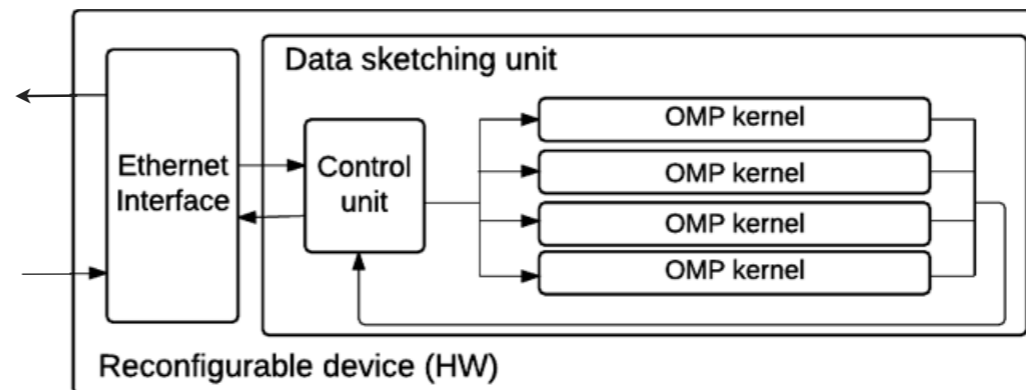
Data Sketching

- Computing the block-sparse matrix V
 - Applying the efficient, and greedy Orthogonal Matching Pursuit (OMP) routine
- Asynchronous parallel approach via a control unit



Hardware Implementation

- Xilinx Virtex-6 FPGA ML605 Evaluation
- IEEE 754 single precision floating point format
- Resource utilization



Virtex-6 resource utilization

	Used	Available	Utilization
Slice Registers	50888	301440	16%
Slice LUTs	81585	150720	54%
RAM B36E1	382	416	91%
DSP 48E1s	356	768	46%

OMP Routine

- OMP is mainly consists of three steps:
 - Find best fitting column
 - LS optimization
 - Residual update
- Use QR decomposition to address the LS optimization

Algorithm 2 OMP algorithm

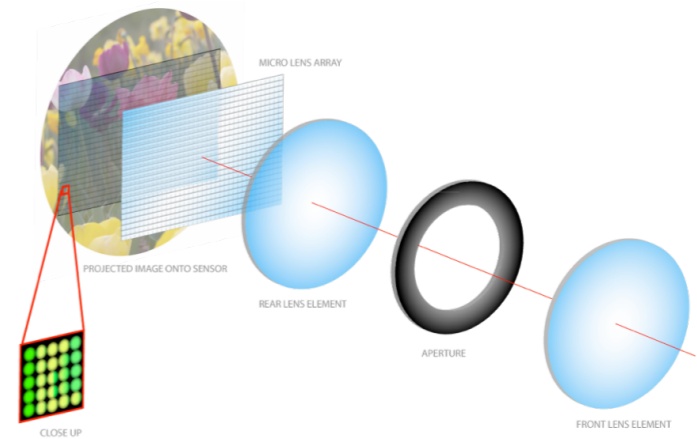
Inputs: Matrix \mathbf{D} , measurement \mathbf{A}_i , sparsity level k , threshold error ϵ .

Output: Support set Λ and k -dimensional coefficient vector \mathbf{v} .

```
1:  $r \leftarrow \mathbf{A}_i$ 
2:  $\Lambda^0 \leftarrow \emptyset$ 
3: for  $i = 1, \dots, k$  do
4:    $\Lambda \leftarrow \Lambda \cup \operatorname{argmax}_j | \langle r^{i-1}, \mathbf{D}_j \rangle |$  Find best fitting column
5:    $v^i \leftarrow \operatorname{argmin}_v \| r^{i-1} - \mathbf{D}_{\Lambda^i} v \|_2^2$  LS Optimization
6:    $r^i \leftarrow r^{i-1} - \mathbf{D}_{\Lambda^i} v^i$  Residual Update
end for
```

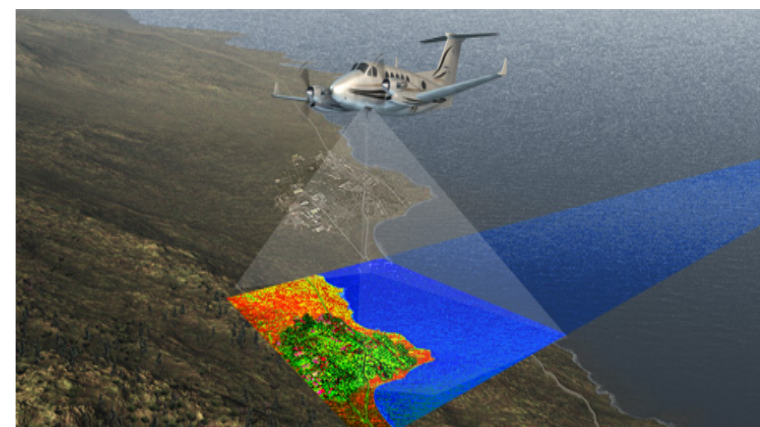
Benchmark Datasets

- Three different datasets:
 - Light Field imaging
 - A sequence of multi-dimensional array of images that are simultaneously captured from slightly different viewpoints



- Hyper-Spectral imaging
 - A sequence of images generated by hundreds of detectors that capture the information from across the electromagnetic spectrum

- Synthetic data



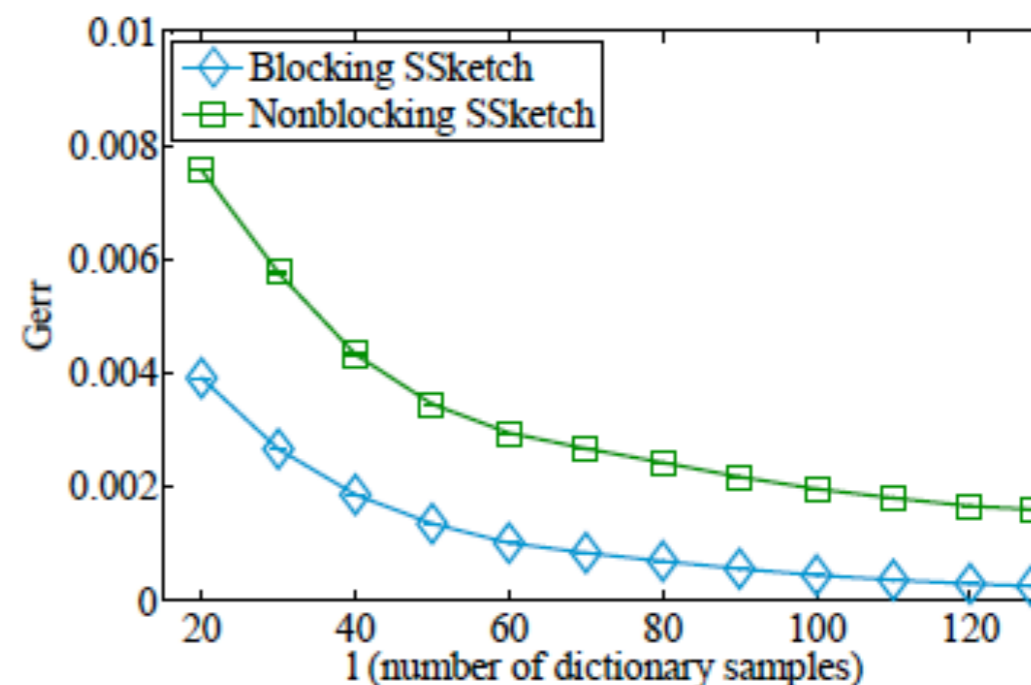
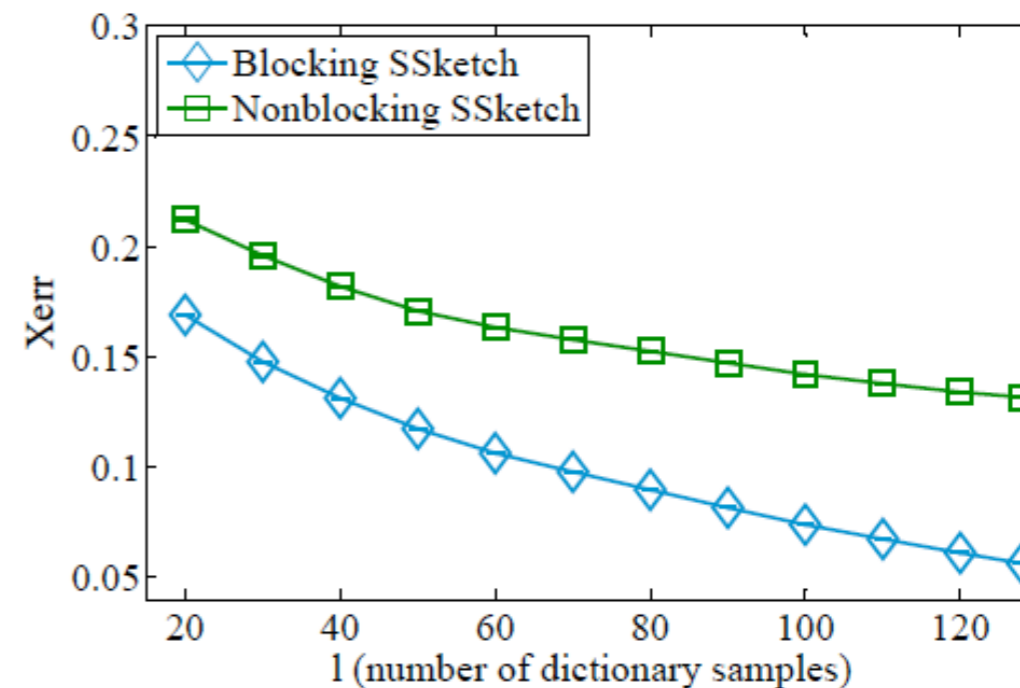
Evaluation Results

$$Xerr = \frac{\|\mathbf{A} - \tilde{\mathbf{A}}\|_F}{\|\mathbf{A}\|_F}, \quad \tilde{\mathbf{A}} = \mathbf{D}\mathbf{V}$$

$$Gerr = \frac{\|\mathbf{A}^t \mathbf{A} - \tilde{\mathbf{A}}^T \tilde{\mathbf{A}}\|_F}{\|\mathbf{A}^t \mathbf{A}\|_F}$$

$$Compression_rate = \frac{nnz(\mathbf{D}) + nnz(\mathbf{V})}{nnz(\mathbf{A})}$$

- Data sketching error decreases as the dictionary size increases



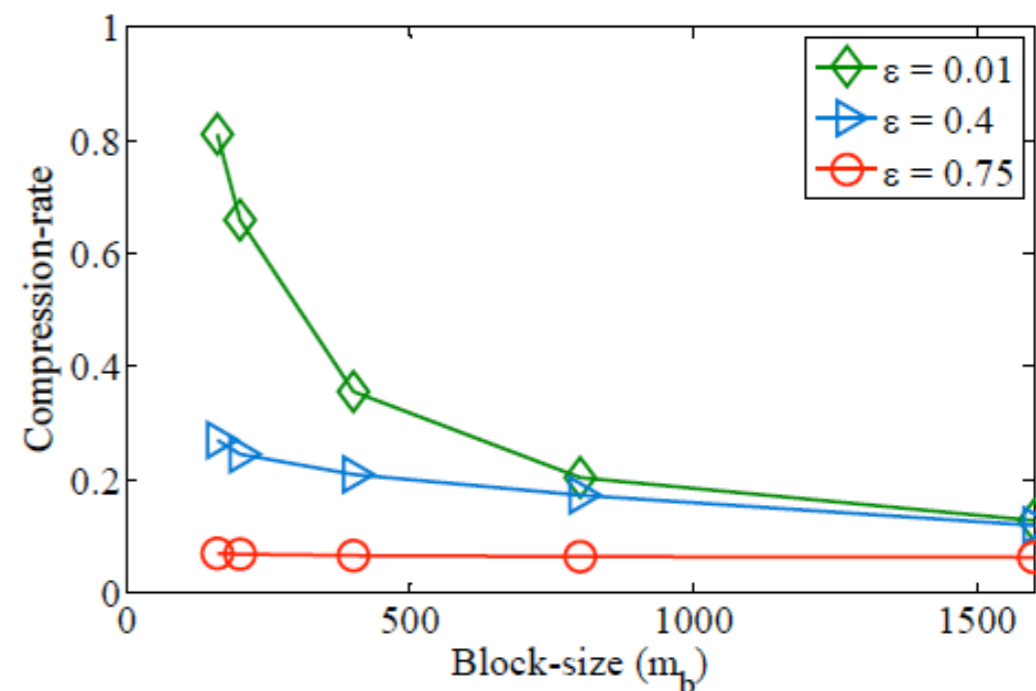
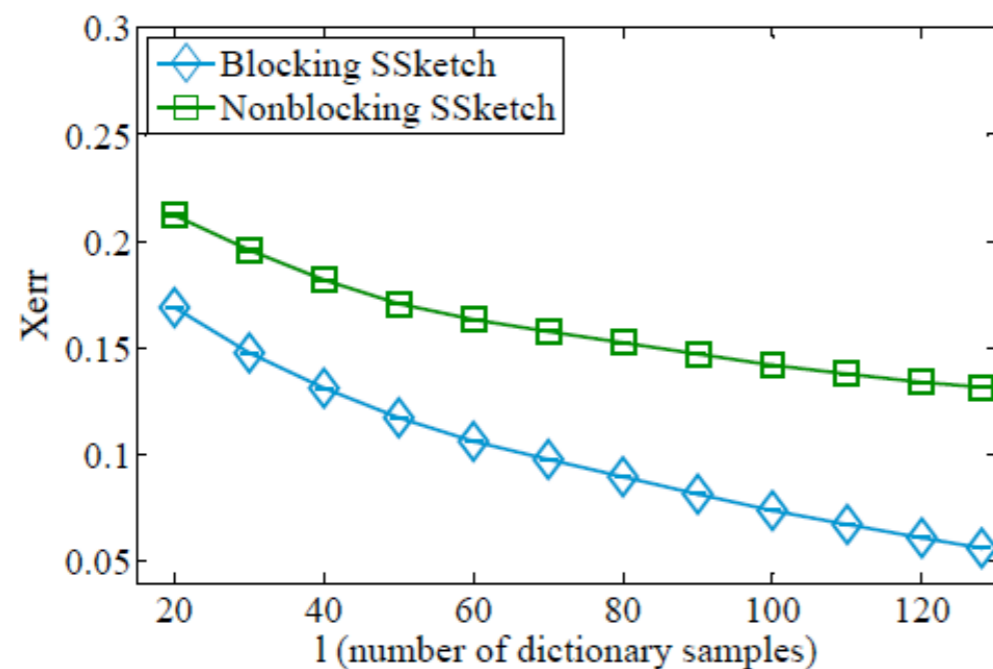
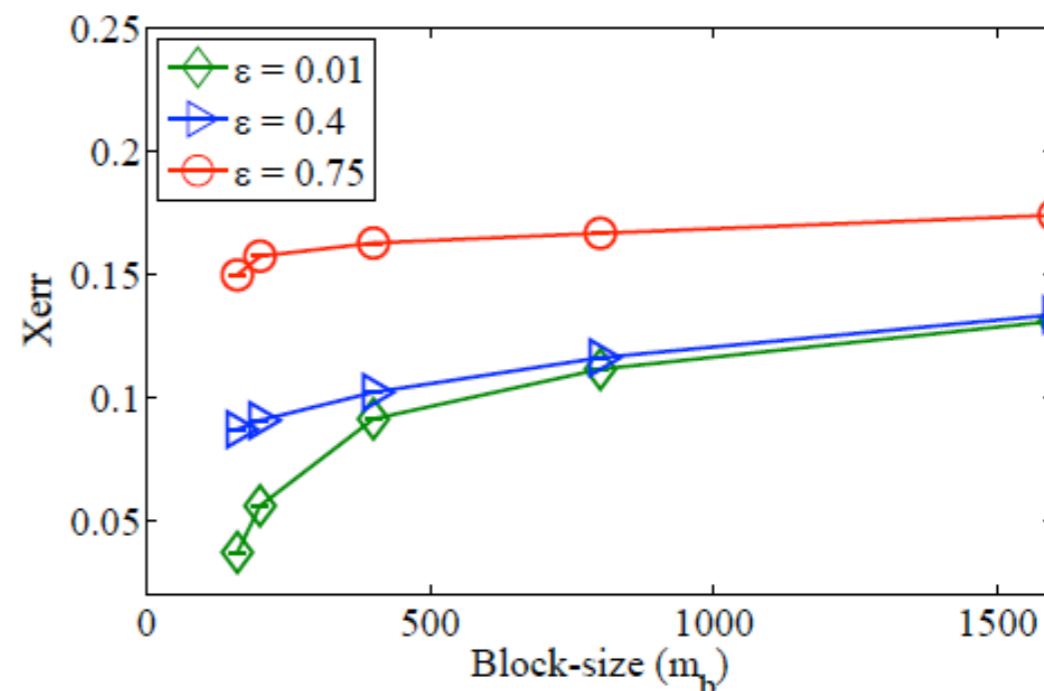
Evaluation Results

$$X_{err} = \frac{\|\mathbf{A} - \tilde{\mathbf{A}}\|_F}{\|\mathbf{A}\|_F}, \quad \tilde{\mathbf{A}} = \mathbf{D}\mathbf{V}$$

$$G_{err} = \frac{\|\mathbf{A}^t \mathbf{A} - \tilde{\mathbf{A}}^T \tilde{\mathbf{A}}\|_F}{\|\mathbf{A}^t \mathbf{A}\|_F}$$

$$Compression_rate = \frac{nnz(\mathbf{D}) + nnz(\mathbf{V})}{nnz(\mathbf{A})}$$

- There is a trade-off between the sketch accuracy and the number of non-zeros in the block-sparse matrix \mathbf{V}



Evaluation Results

- SSketch total processing time is linear in terms of the number of processed samples
- Runtime: $T_{SSketch} \approx T_{\substack{\text{dictionary} \\ \text{learning}}} + T_{\substack{\text{Communication} \\ \text{Overhead}}} + T_{\substack{\text{FPGA} \\ \text{Computation}}}$
- SSketch runtime is dominated by $T_{\substack{\text{FPGA} \\ \text{Computation}}}$

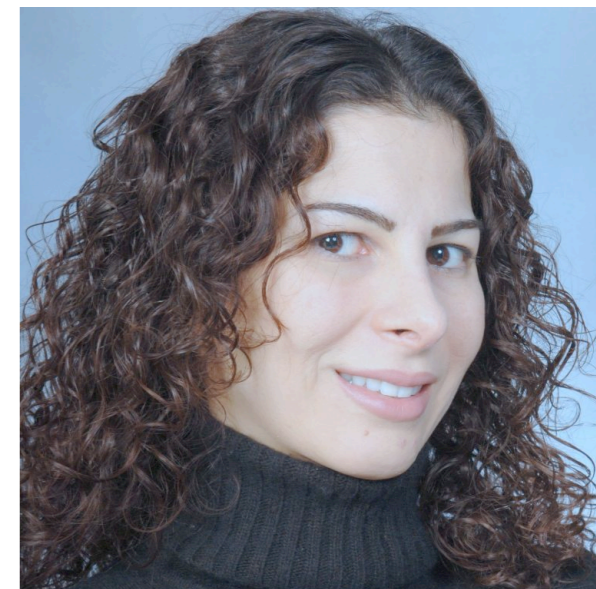
<i>Size of n</i>	$T_{SSketch}$ ($l = 128$)	$T_{SSketch}$ ($l = 64$)
$1k$	3.635s	2.31s
$5k$	21.029s	12.01s
$10k$	43.446s	24.32s
$20k$	90.761s	48.52s

$m_b = 256, \epsilon = 0.01, \text{ and } \alpha = 0.1$

Conclusion

- Adaptive hardware-accelerated streaming-based data transformation
- Scalable streaming-based sketching methodology
 - Amenable to FPGA acceleration
 - Fixed, and low memory footprint
- User-friendly API
 - Rapid prototyping of an arbitrary matrix-based data analysis
- Scalable, floating-point implementation of OMP on FPGA
- Up to 200 folds speed up compared to the software-only realization
- Less than 4% message passing delay for communication between the processor and accelerator

Acknowledgment



- From left to right:

Bita Darvish Rouhani, PhD. student, Rice University (bita@rice.edu)

Ebrahim Songhori, PhD. student, Rice University (ebrahim@rice.edu)

Azalia Mirhoseini, PhD. student, Rice University (azalia@rice.edu)

Farinaz Koushanfar, Associate Professor, Rice University (farinaz@rice.edu)

- This work was supported in parts by the Office of Naval Research grant (ONR- N00014-11-1-0885)