

Space Efficient TREC for Enabling Deep Learning on Microcontrollers

Jiawei Guan \diamond , Feng Zhang \diamond , Jiesong Liu \diamond , Hsin-Hsuan Sung \star ,
Ruofan Wu \diamond , Xiaoyong Du \diamond , Xipeng Shen \star

\diamond Key Laboratory of Data Engineering and Knowledge Engineering (MOE),
and School of Information, Renmin University of China
 \star Computer Science Department, North Carolina State University



Outline

- Motivation
- TREC Architecture
- Experiments
- Results

Motivation

AIoT

- CNN **Inference** on the edge.



Motivation

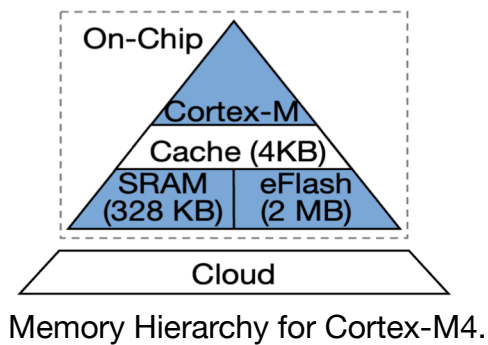
AIoT

- CNN **Inference** on the edge.
- Computation-intensive.

Motivation

AIoT

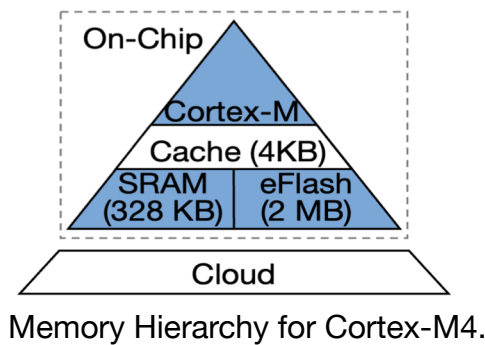
- CNN **Inference** on the edge.
- Computation-intensive.
- Resource-constrained devices.



Motivation

AIoT

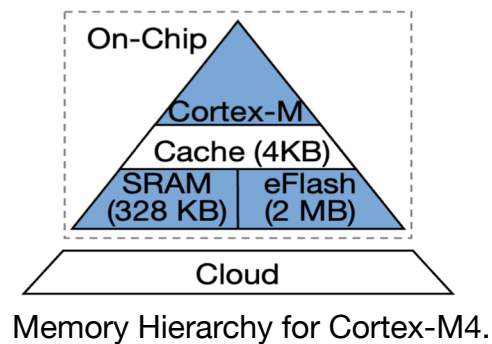
- CNN **Inference** on the edge.
- Computation-intensive.
- Resource-constrained devices.
- Large models are required to be compressed.



Motivation

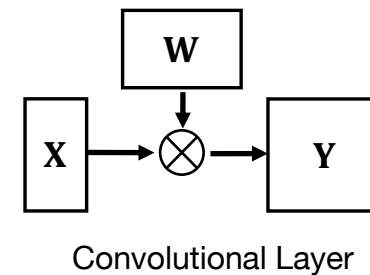
AIoT

- CNN **Inference** on the edge.
- Computation-intensive.
- Resource-constrained devices.
- Large models are required to be compressed.



Redundancy

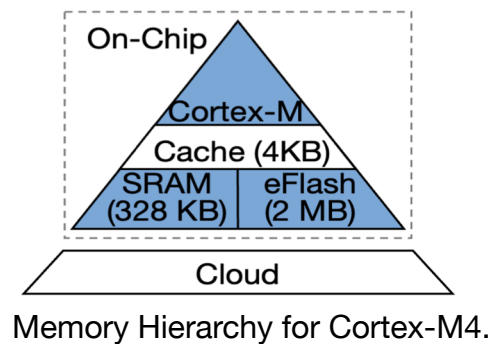
- Lasting redundancy
- Transient redundancy



Motivation

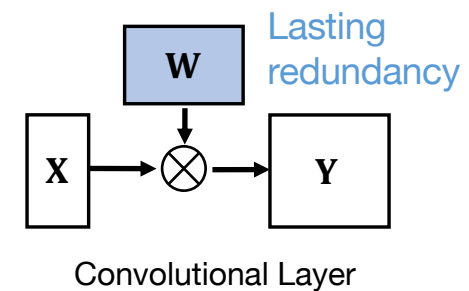
AIoT

- CNN **Inference** on the edge.
- Computation-intensive.
- Resource-constrained devices.
- Large models are required to be compressed.



Redundancy

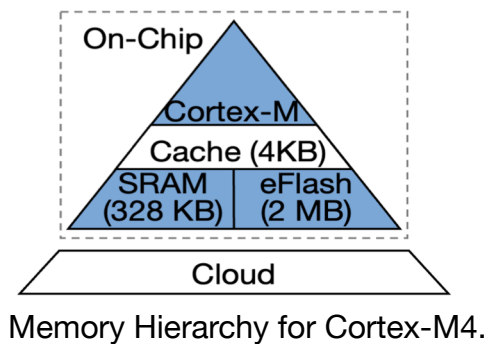
- Lasting redundancy
 - Arises from DNN **parameters**.
 - Removed during training mostly.
 - Pruning, quantization, etc.
- Transient redundancy



Motivation

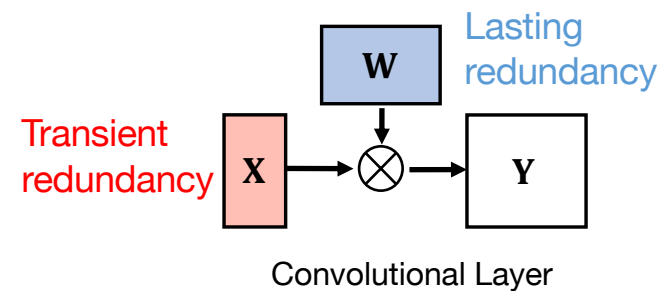
AIoT

- CNN **Inference** on the edge.
- Computation-intensive.
- Resource-constrained devices.
- Large models are required to be compressed.



Redundancy

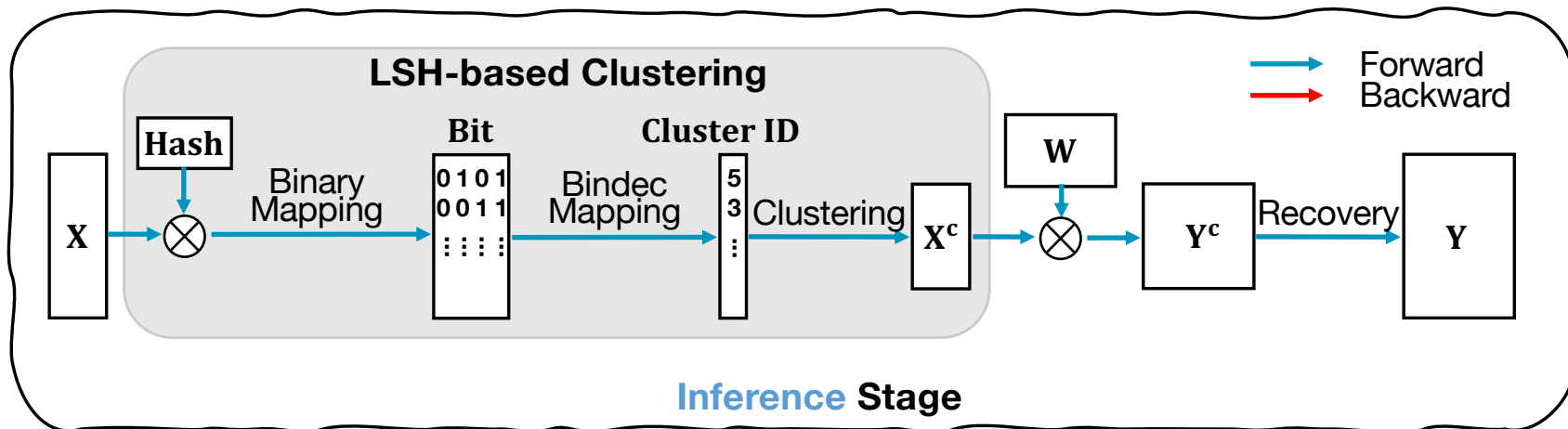
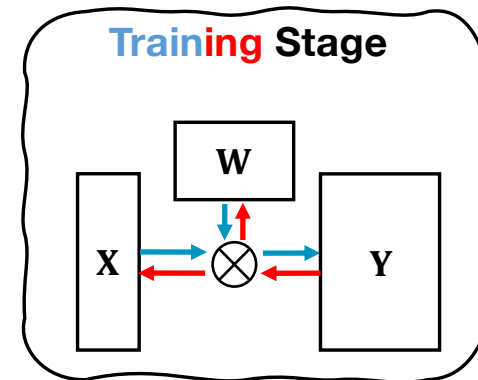
- Lasting redundancy
 - Arises from DNN **parameters**.
 - Removed during training mostly.
 - Pruning, quantization, etc.
- Transient redundancy
 - Arises from **Input data / activation maps**.
 - Removed during inference.



Motivation

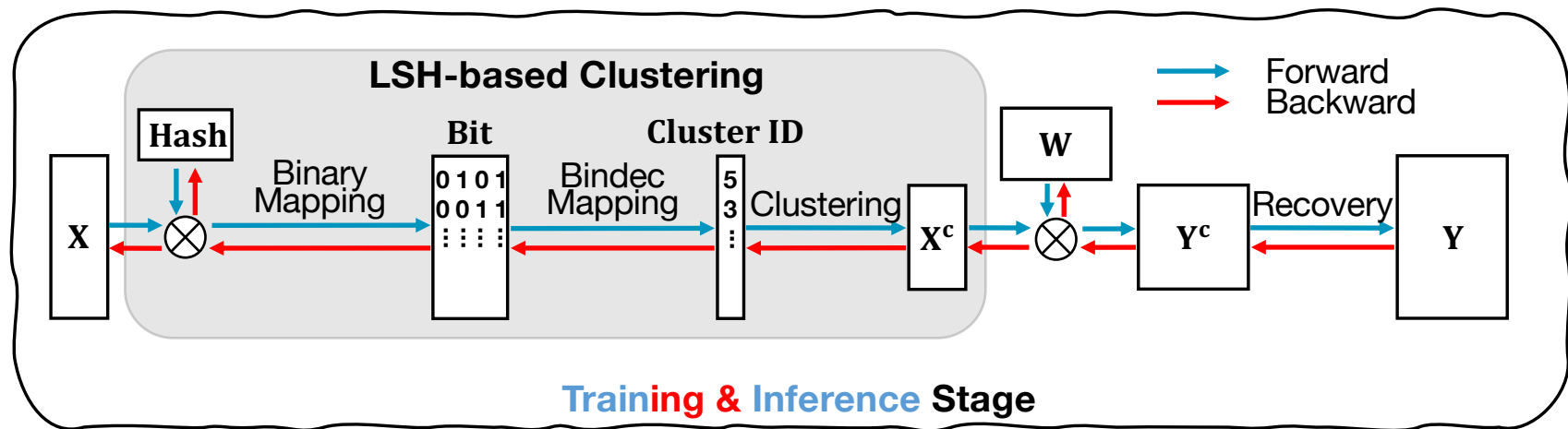
DeepReuse

- Treat transient redundancy in an **Ad-hoc** manner.
- Over 5% **accuracy fluctuations** in different runs.

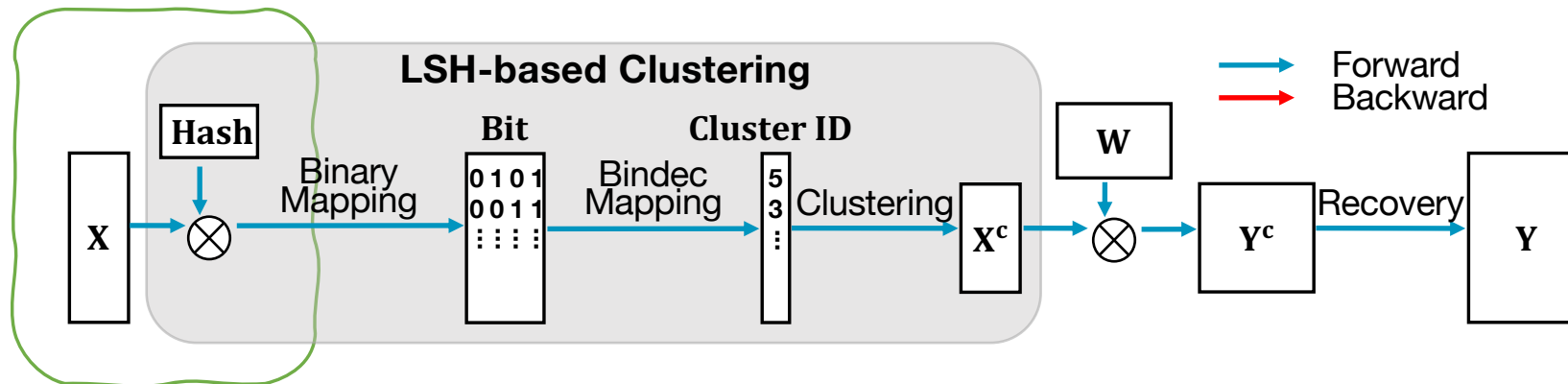


Concept

Can the detection and elimination of transient redundancy be integrated into the CNN training process, so that they become part of its inherent architecture?

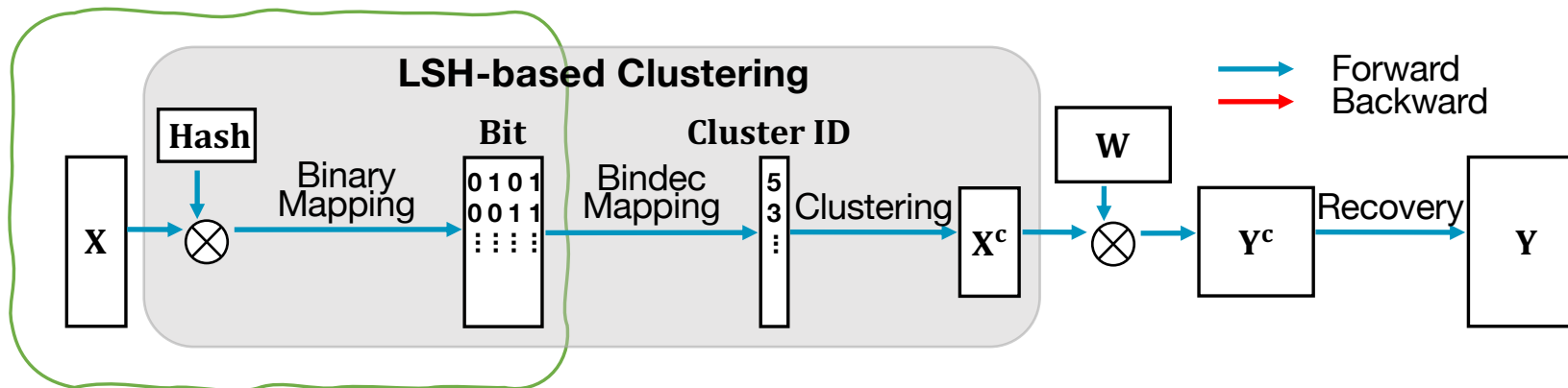


A Closer Look at Redundancy Elimination



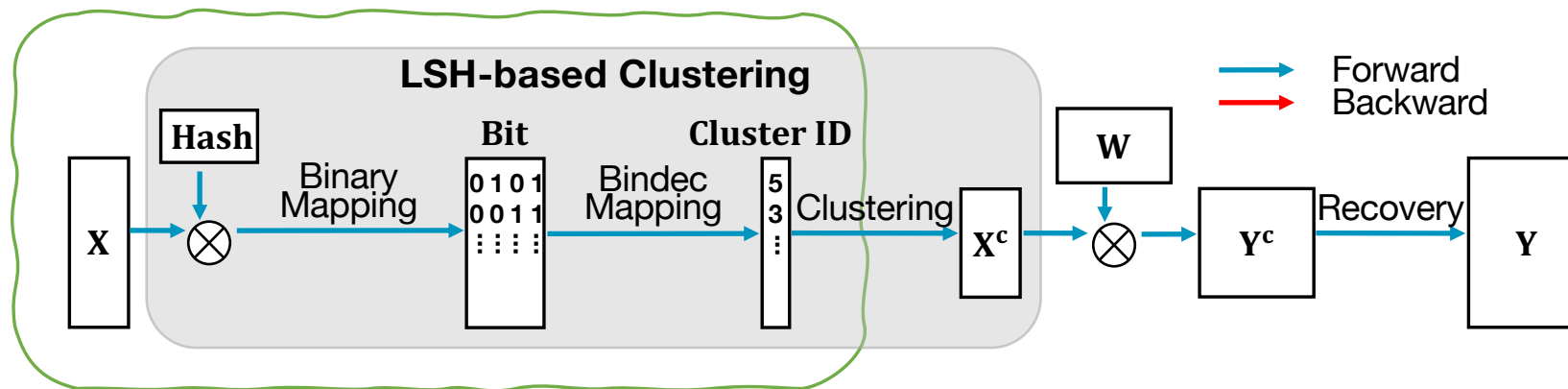
$$\begin{matrix}
 \mathbf{X} \\
 \begin{bmatrix} 1 & 2 & 3 \\ 0 & 0 & 1 \\ 0 & 2 & 1 \\ 1 & 3 & 2 \end{bmatrix}
 \end{matrix}
 \times
 \begin{matrix}
 \text{Hash} \\
 \begin{bmatrix} -1 & 1 \\ 1 & 2 \\ 0 & -1 \end{bmatrix}
 \end{matrix}
 =
 \begin{matrix}
 \text{Projected} \\
 \begin{bmatrix} 1 & 2 \\ 0 & -1 \\ 2 & 3 \\ 2 & 5 \end{bmatrix}
 \end{matrix}$$

A Closer Look at Redundancy Elimination



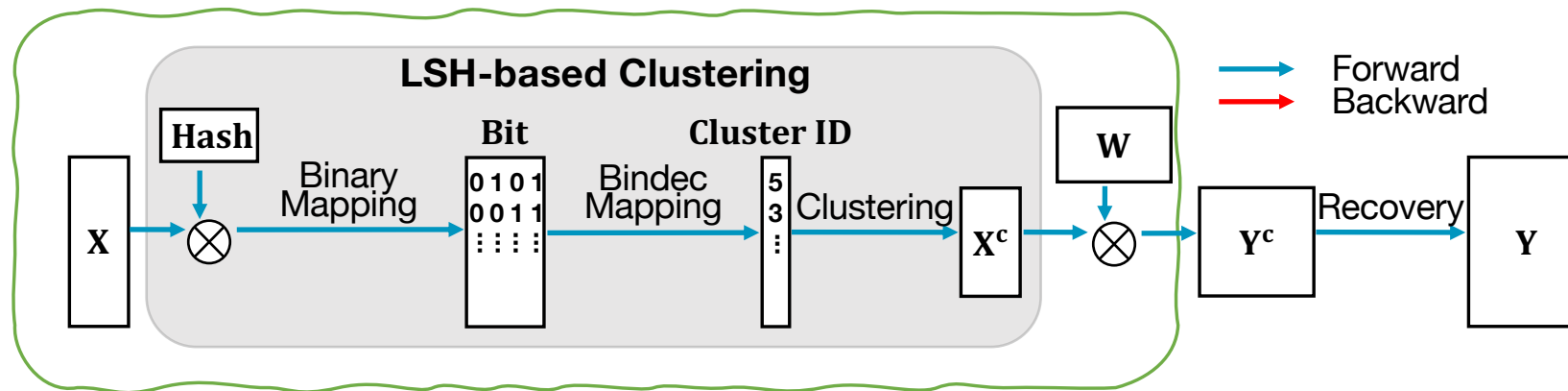
$$\begin{array}{c}
 \mathbf{X} \\
 \begin{bmatrix} 1 & 2 & 3 \\ 0 & 0 & 1 \\ 0 & 2 & 1 \\ 1 & 3 & 2 \end{bmatrix}
 \end{array}
 \times
 \begin{array}{c}
 \text{Hash} \\
 \begin{bmatrix} -1 & 1 \\ 1 & 2 \\ 0 & -1 \end{bmatrix}
 \end{array}
 =
 \begin{array}{c}
 \text{Projected} \\
 \begin{bmatrix} 1 & 2 \\ 0 & -1 \\ 2 & 3 \\ 2 & 5 \end{bmatrix}
 \end{array}
 \xrightarrow{\text{Binary Mapping}}
 \begin{array}{c}
 \text{Bit} \\
 \begin{bmatrix} 1 & 1 \\ 0 & 0 \\ 1 & 1 \\ 1 & 1 \end{bmatrix}
 \end{array}$$

A Closer Look at Redundancy Elimination



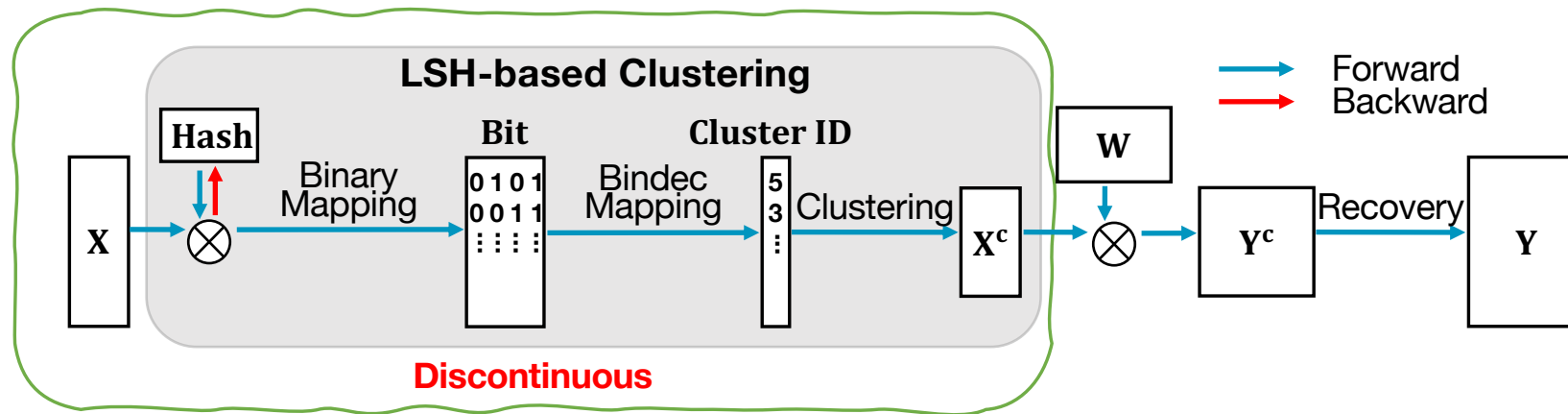
$$\begin{array}{c}
 \mathbf{X} \\
 \begin{bmatrix} 1 & 2 & 3 \\ 0 & 0 & 1 \\ 0 & 2 & 1 \\ 1 & 3 & 2 \end{bmatrix}
 \end{array}
 \times
 \begin{array}{c}
 \text{Hash} \\
 \begin{bmatrix} -1 & 1 \\ 1 & 2 \\ 0 & -1 \end{bmatrix}
 \end{array}
 =
 \begin{array}{c}
 \text{Projected} \\
 \begin{bmatrix} 1 & 2 \\ 0 & -1 \\ 2 & 3 \\ 2 & 5 \end{bmatrix}
 \end{array}
 \xrightarrow{\text{Binary Mapping}}
 \begin{array}{c}
 \text{Bit} \\
 \begin{bmatrix} 1 & 1 \\ 0 & 0 \\ 1 & 1 \\ 1 & 1 \end{bmatrix}
 \end{array}
 \xrightarrow{\text{Bindec Mapping}}
 \begin{array}{c}
 \text{Cluster ID} \\
 \begin{bmatrix} 3 \\ 0 \\ 3 \\ 3 \end{bmatrix}
 \end{array}$$

A Closer Look at Redundancy Elimination

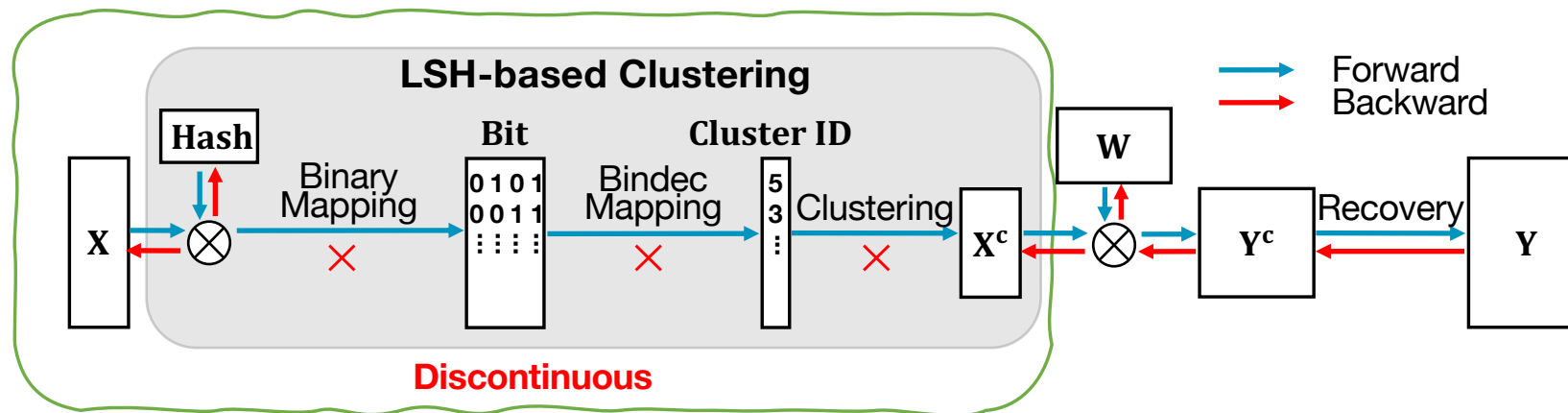


$$\begin{array}{c}
 \mathbf{X} \\
 \begin{bmatrix} 1 & 2 & 3 \\ 0 & 0 & 1 \\ 0 & 2 & 1 \\ 1 & 3 & 2 \end{bmatrix}
 \end{array}
 \times
 \begin{array}{c}
 \text{Hash} \\
 \begin{bmatrix} -1 & 1 \\ 1 & 2 \\ 0 & -1 \end{bmatrix}
 \end{array}
 =
 \begin{array}{c}
 \text{Projected} \\
 \begin{bmatrix} 1 & 2 \\ 0 & -1 \\ 2 & 3 \\ 2 & 5 \end{bmatrix}
 \end{array}
 \xrightarrow{\text{Binary Mapping}}
 \begin{array}{c}
 \text{Bit} \\
 \begin{bmatrix} 1 & 1 \\ 0 & 0 \\ 1 & 1 \\ 1 & 1 \end{bmatrix}
 \end{array}
 \xrightarrow{\text{Bindex Mapping}}
 \begin{array}{c}
 \text{Cluster ID} \\
 \begin{bmatrix} 3 \\ 0 \\ 3 \\ 3 \end{bmatrix}
 \end{array}$$

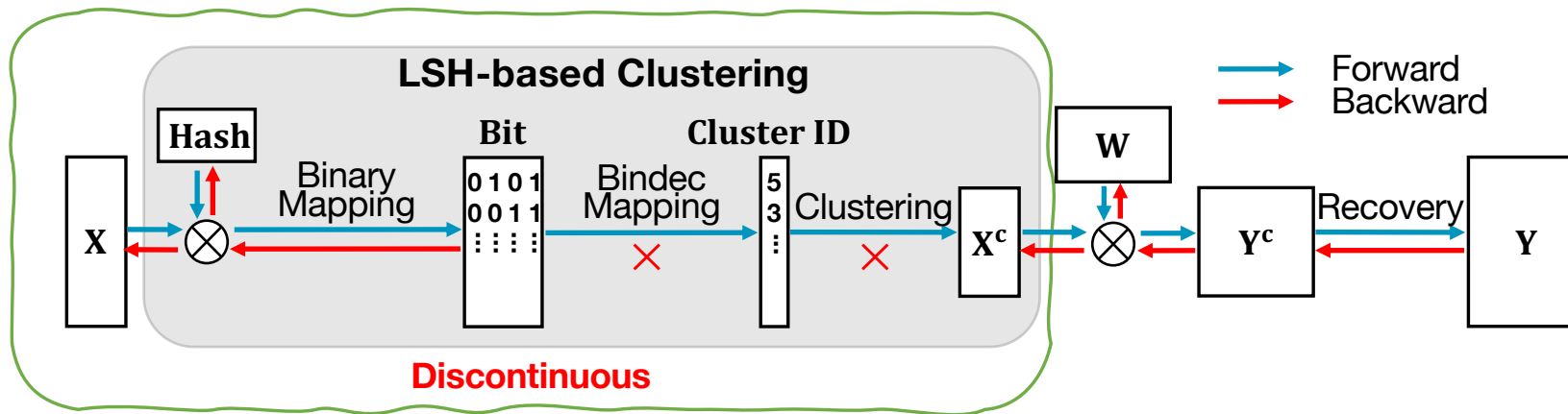
Backward Pass



Backward Pass

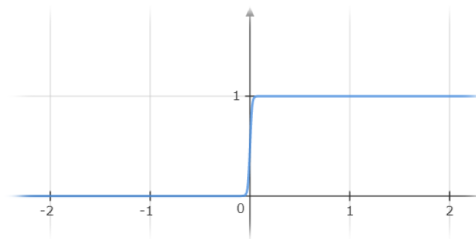


Backward Pass

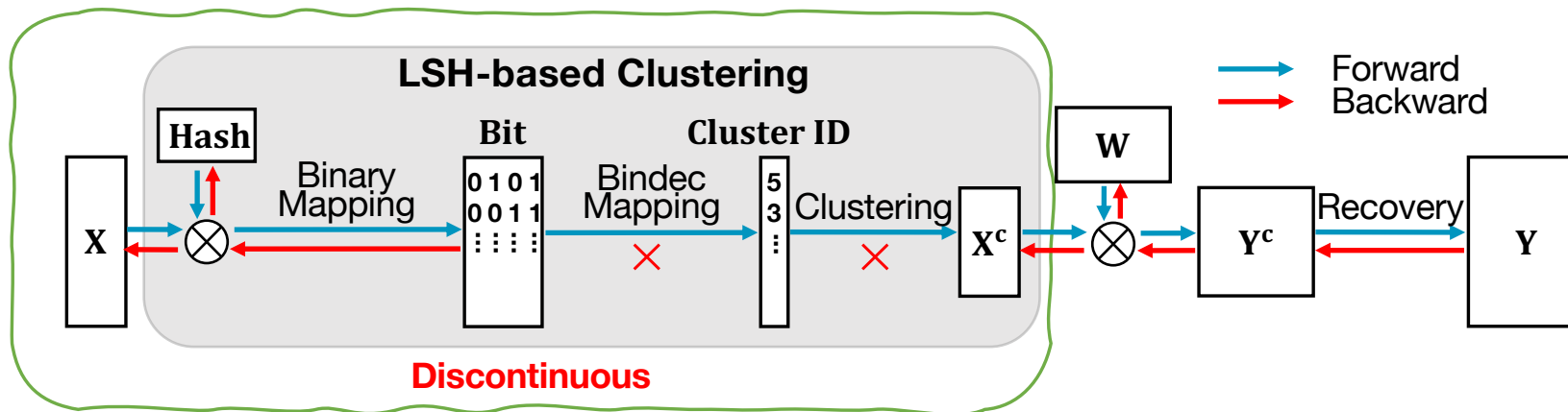


1. Binary Mapping \rightarrow Binary Approximation

$$\text{sigmoid}(x) = \frac{1}{1 + e^{-\alpha x}}$$



Backward Pass

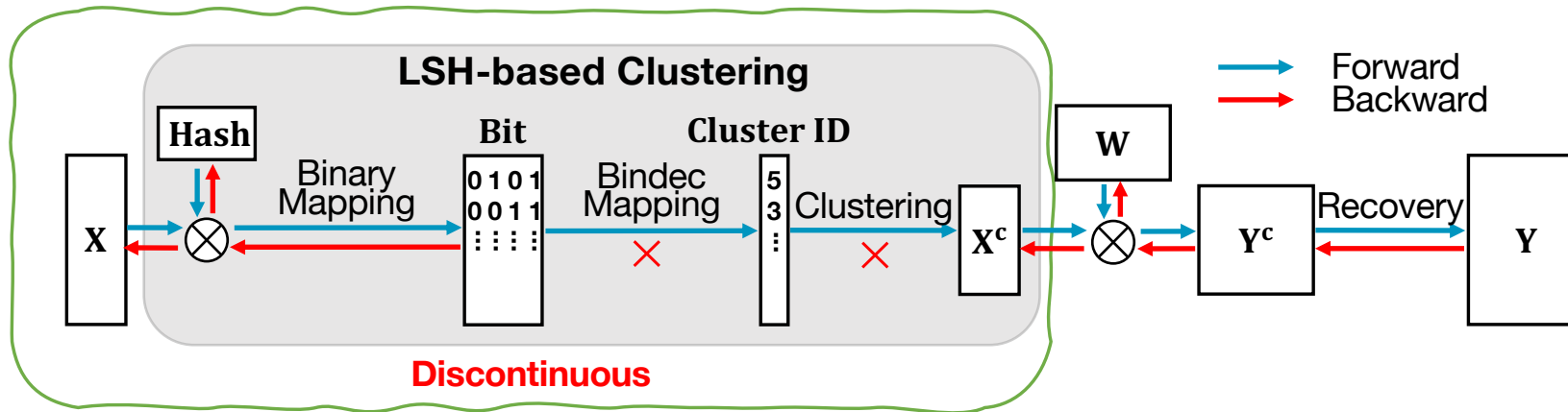


2. Bindex Mapping \rightarrow Bindex Conversion

- Transformation matrix

$$\begin{bmatrix} 2^{H-1}/1 & 2^{H-1}/2 & \dots & 2^{H-1}/2^H \\ \vdots & \vdots & \ddots & \vdots \\ 2^1/1 & 2^1/2 & \dots & 2^1/2^H \\ 2^0/1 & 2^0/2 & \dots & 2^0/2^H \end{bmatrix}$$

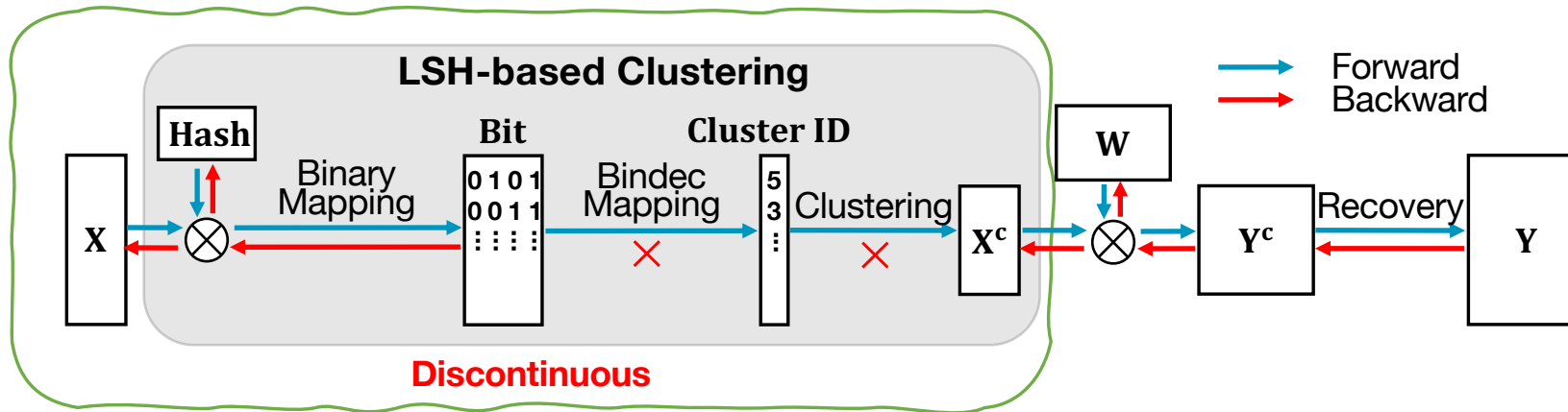
Backward Pass



Bindec Conversion

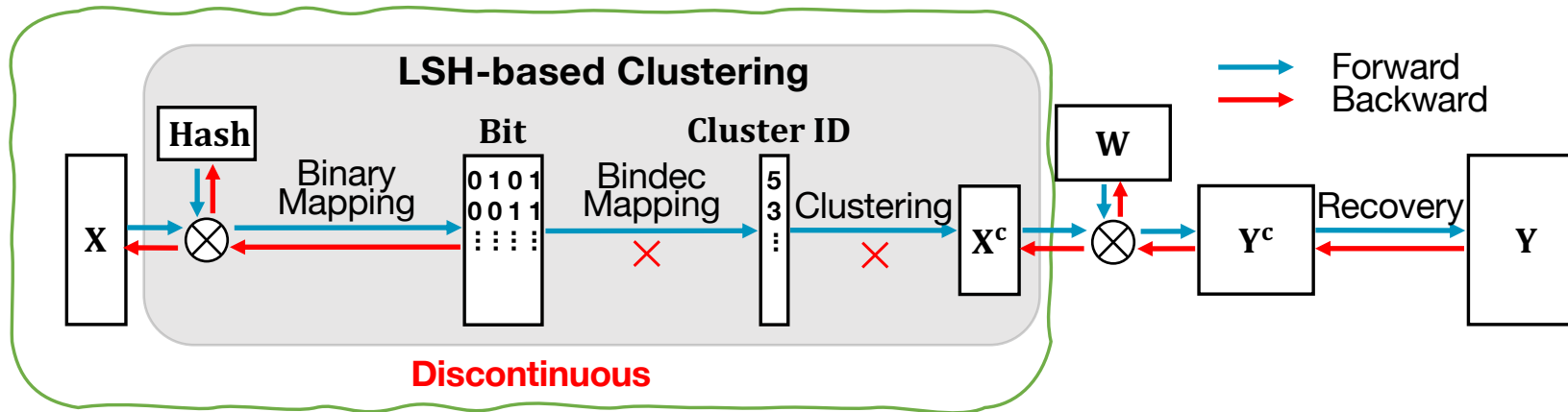
$$\begin{array}{c} \text{Bit} \\ \left[\begin{array}{cc} 1 & 1 \\ 0 & 0 \\ 1 & 1 \\ 1 & 1 \end{array} \right] + \left[\begin{array}{cc} 0 & 1 \\ 0 & 1 \\ 0 & 1 \\ 0 & 1 \end{array} \right] \end{array} \times \begin{array}{c} \text{Transformation} \\ (\mathbf{H} \times 2^H) \\ \left[\begin{array}{cccc} 2^1/1 & 2^1/2 & 2^1/3 & 2^1/4 \\ 2^0/1 & 2^0/2 & 2^0/3 & 2^0/4 \end{array} \right] \end{array} = \begin{array}{c} \text{Quotient} \\ \left[\begin{array}{cccc} 4 & 2 & 4/3 & 1 \\ 1 & 1/2 & 1/3 & 1/4 \\ 4 & 2 & 4/3 & 1 \\ 4 & 2 & 4/3 & 1 \end{array} \right] \end{array}$$

Backward Pass



$$\begin{array}{c} \text{Bit} \\ \left[\begin{array}{cc} 1 & 1 \\ 0 & 0 \\ 1 & 1 \\ 1 & 1 \end{array} \right] + \left[\begin{array}{cc} 0 & 1 \\ 0 & 1 \\ 0 & 1 \\ 0 & 1 \end{array} \right] \end{array} \times \begin{array}{c} \text{Transformation} \\ (H \times 2^H) \\ \left[\begin{array}{cc|cc} 2^1/ & 2^1/ & 2^1/ & 2^1/ \\ \hline 2^0/ & 2^0/ & 2^0/ & 2^0/ \end{array} \right] \end{array} = \begin{array}{c} \text{Quotient} \\ \left[\begin{array}{cc|cc} 4 & 4 & 4 & 4 \\ \hline 1 & 1 & 1 & 1 \\ \hline 4 & 4 & 4 & 4 \\ \hline 4 & 4 & 4 & 4 \end{array} \right] \end{array}$$

Backward Pass



Bit

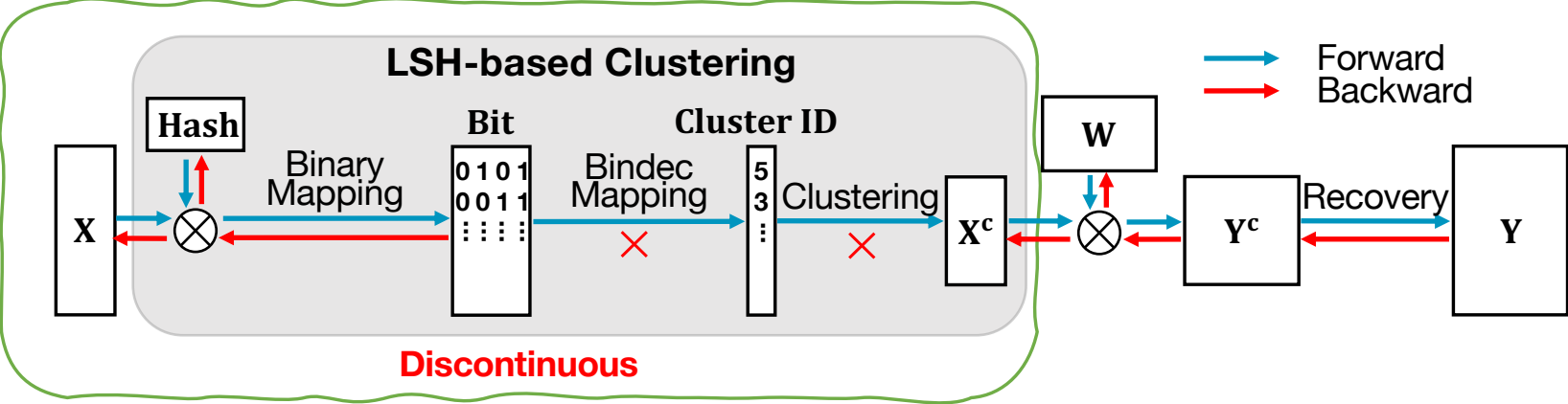
$$\begin{bmatrix} 1 & 1 \\ 0 & 0 \\ 1 & 1 \\ 1 & 1 \end{bmatrix} + \begin{bmatrix} 0 & 1 \\ 0 & 1 \\ 0 & 1 \\ 0 & 1 \end{bmatrix} \times \begin{bmatrix} 2^1/ & 2^1/ & 2^1/ & 2^1/ \\ 2^0/ & 2^0/ & 2^0/ & 2^0/ \end{bmatrix} = \begin{bmatrix} 4 & 4 & 4 & 4 \\ 1 & 1 & 1 & 1 \\ 4 & 4 & 4 & 4 \\ 4 & 4 & 4 & 4 \end{bmatrix}$$

Transformation
($H \times 2^H$)

Quotient

$$\begin{bmatrix} 4 & 4 & 4 & 4 \\ 1 & 1 & 1 & 1 \\ 4 & 4 & 4 & 4 \\ 4 & 4 & 4 & 4 \end{bmatrix} \times \begin{bmatrix} 1/1 & 1/2 & 1/3 & 1/4 \\ 1/1 & 1/2 & 1/3 & 1/4 \end{bmatrix} = \begin{bmatrix} 4 & 2 & 4/3 & 1 \\ 1 & 1/2 & 1/3 & 1/4 \\ 4 & 2 & 4/3 & 1 \\ 4 & 2 & 4/3 & 1 \end{bmatrix}$$

Backward Pass



Bit

$$\begin{bmatrix} 1 & 1 \\ 0 & 0 \\ 1 & 1 \\ 1 & 1 \end{bmatrix} + \begin{bmatrix} 0 & 1 \\ 0 & 1 \\ 0 & 1 \\ 0 & 1 \end{bmatrix}$$

Transformation
($H \times 2^H$)

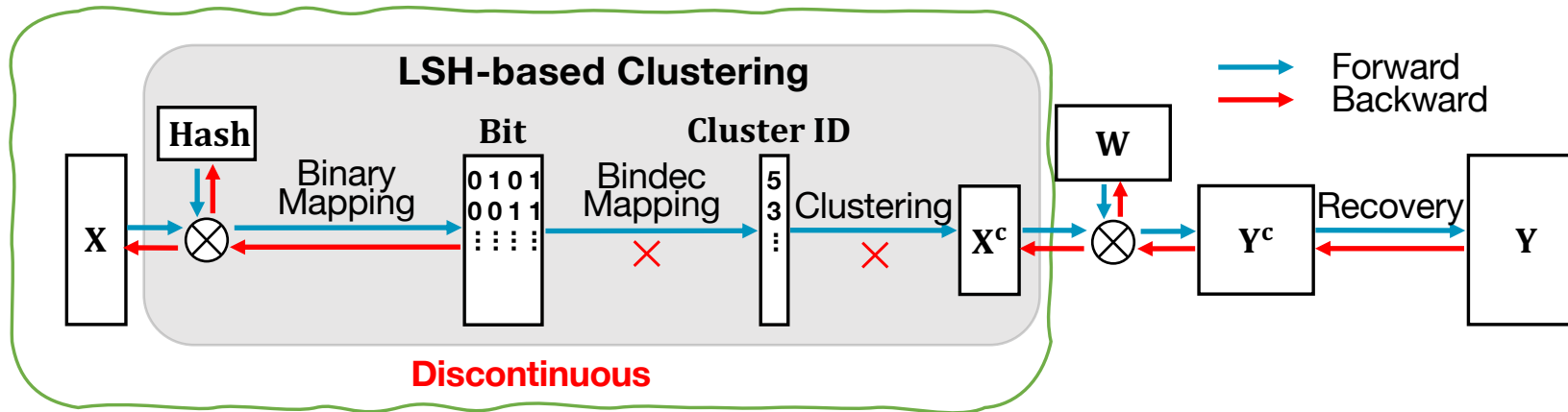
$$\times \begin{bmatrix} 2^1/1 & 2^1/2 & 2^1/3 & 2^1/4 \\ 2^0/1 & 2^0/2 & 2^0/3 & 2^0/4 \end{bmatrix}$$

Space waste

Quotient

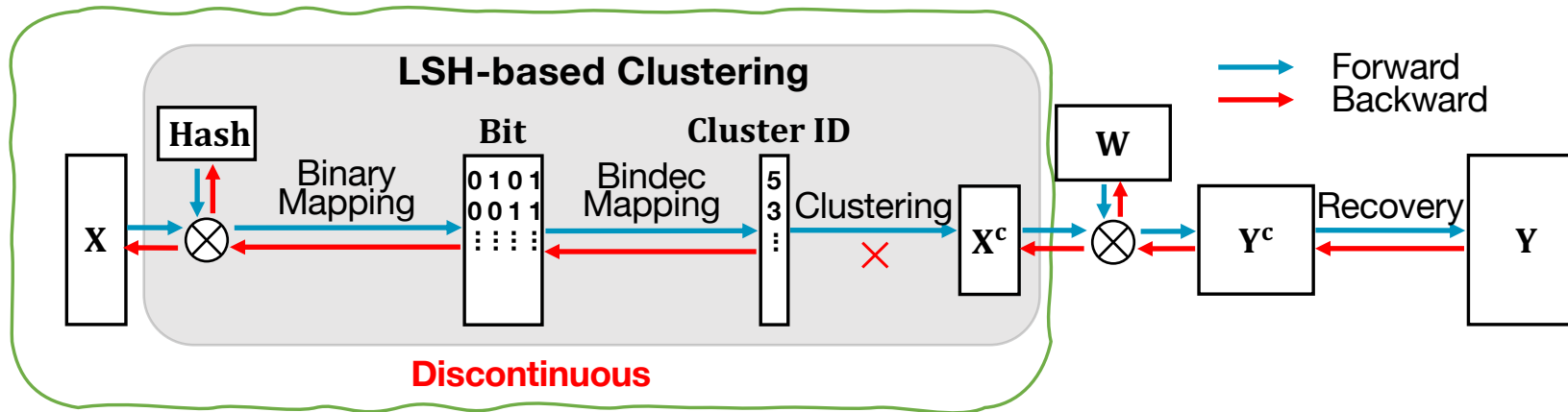
$$= \begin{bmatrix} 4 & 2 & 4/3 & 1 \\ 1 & 1/2 & 1/3 & 1/4 \\ 4 & 2 & 4/3 & 1 \\ 4 & 2 & 4/3 & 1 \end{bmatrix}$$

Backward Pass



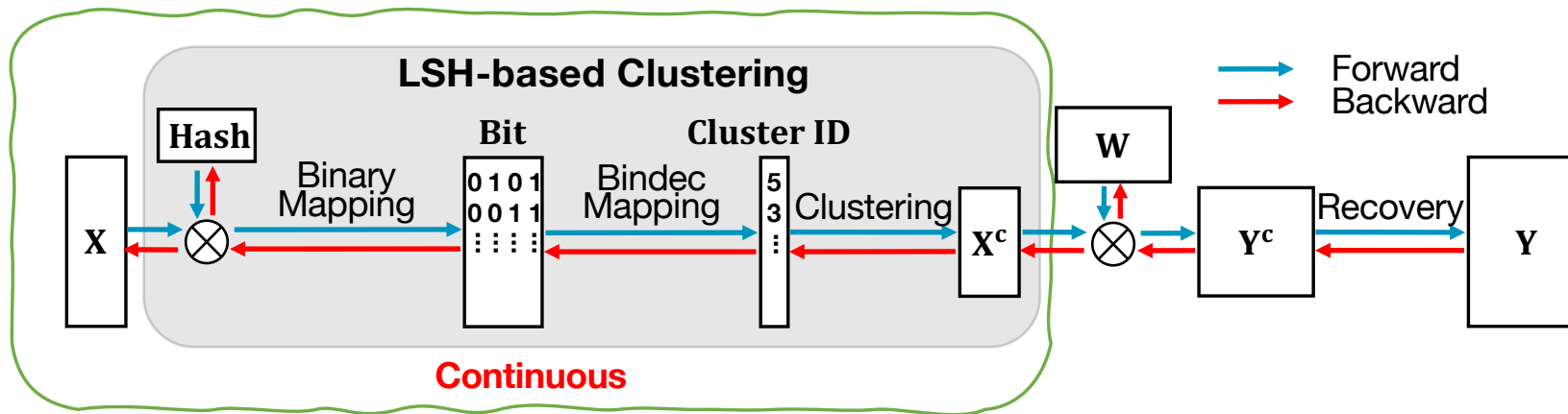
$$\begin{array}{c}
 \text{Bit} \\
 \left[\begin{array}{cc} 1 & 1 \\ 0 & 0 \\ 1 & 1 \\ 1 & 1 \end{array} \right] + \left[\begin{array}{cc} 0 & 1 \\ 0 & 1 \\ 0 & 1 \\ 0 & 1 \end{array} \right] \\
 \times \left[\begin{array}{cccc} 2^1/1 & 2^1/2 & 2^1/3 & 2^1/4 \\ 2^0/1 & 2^0/2 & 2^0/3 & 2^0/4 \end{array} \right] = \left[\begin{array}{cccc} 4 & 2 & 4/3 & 1 \\ 1 & 1/2 & 1/3 & 1/4 \\ 4 & 2 & 4/3 & 1 \\ 4 & 2 & 4/3 & 1 \end{array} \right] \\
 \downarrow \text{Forward information} \\
 \left[\begin{array}{cc} 1 & 1 \\ 0 & 0 \\ 1 & 1 \\ 1 & 1 \end{array} \right] + \left[\begin{array}{cc} 0 & 1 \\ 0 & 1 \\ 0 & 1 \\ 0 & 1 \end{array} \right] \times \left[\begin{array}{cc} 2^1/1 & 2^1/4 \\ 2^0/1 & 2^0/4 \end{array} \right] = \left[\begin{array}{cc} 4 & 1 \\ 1 & 1/4 \\ 4 & 1 \\ 4 & 1 \end{array} \right]
 \end{array}$$

Backward Pass



$$\begin{aligned}
 & \text{Bit} \quad \begin{bmatrix} 1 & 1 \\ 0 & 0 \\ 1 & 1 \\ 1 & 1 \end{bmatrix} + \begin{bmatrix} 0 & 1 \\ 0 & 1 \\ 0 & 1 \\ 0 & 1 \end{bmatrix} \times \begin{bmatrix} 2^1/1 & 2^1/2 & 2^1/3 & 2^1/4 \\ 2^0/1 & 2^0/2 & 2^0/3 & 2^0/4 \end{bmatrix} = \begin{bmatrix} 4 & 2 & 4/3 & 1 \\ 1 & 1/2 & 1/3 & 1/4 \\ 4 & 2 & 4/3 & 1 \\ 4 & 2 & 4/3 & 1 \end{bmatrix} \\
 & \quad \downarrow \text{Forward information} \\
 & \begin{bmatrix} 1 & 1 \\ 0 & 0 \\ 1 & 1 \\ 1 & 1 \end{bmatrix} + \begin{bmatrix} 0 & 1 \\ 0 & 1 \\ 0 & 1 \\ 0 & 1 \end{bmatrix} \times \begin{bmatrix} 2^1/1 & 2^1/4 \\ 2^0/1 & 2^0/4 \end{bmatrix} = \begin{bmatrix} 4 & 1 \\ 1 & 1/4 \\ 4 & 1 \\ 4 & 1 \end{bmatrix} \xrightarrow{\text{Gaussian}} \begin{bmatrix} 0 & 1 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \end{bmatrix} \xrightarrow{\text{Divide } N_c} \begin{bmatrix} 0 & 1/3 \\ 1 & 0 \\ 0 & 1/3 \\ 0 & 1/3 \end{bmatrix}_{25} \\
 & \quad \text{Average (Avg)}
 \end{aligned}$$

Backward Pass



- Clustering

$$\begin{bmatrix} 1/3 & 0 & 1/3 & 1/3 \\ 0 & 1 & 0 & 0 \end{bmatrix}^{Avg^T} \times \begin{bmatrix} 1 & 2 & 3 \\ 0 & 0 & 1 \\ 0 & 2 & 1 \\ 1 & 3 & 2 \end{bmatrix}^X = \begin{bmatrix} 2/3 & 7/3 & 2 \\ 2 & 1 & 1 \end{bmatrix}^{X^c}$$

Experiments

Platforms

- **Training:** A server machine with an Intel Core i7-12700K processor and an NVIDIA GeForce RTX A6000 GPU.
- **Inference:** An STM32F469NI MCU.

Datasets

- Cifar-10
- ImageNet-64x64

End-to-end Performance Comparison

Network	Conv Method	Average Time per Image (ms)	Top-1 Accuracy (%)
CifarNet	Conventional	217.32	78.2
	Deep Reuse	154.44	73.2 ~ 76.1
	TREC	153.92	76.5
ZfNet	Conventional	3557.32	80.1
	Deep Reuse	814.03	72.5 ~ 76.6
	TREC	814.01	78.9
Vanilla SqueezeNet	Conventional	1639.51	83.5
	Deep Reuse	328.97	79.8 ~ 81.9
	TREC	327.90	83.0
SqueezeNet + Complex Bypass	Conventional	1998.86	85.3
	Deep Reuse	543.71	80.5 ~ 83.1
	TREC	544.03	84.6
ResNet-34 (ImageNet-64×64)	Conventional	4242.77	52.6
	Deep Reuse	1379.26	46.7 ~ 49.9
	TREC	1378.75	52.2

Experiments

Benefits

- Speedup with minor accuracy loss.
- Stable performance and high accuracy.
- A plug-and-play replacement for convolutional layers in mainstream CNNs.
- Orthogonal to lasting redundancy elimination methods (e.g. pruning, quantization).

End-to-end Performance Comparison

Network	Conv Method	Average Time per Image (ms)	Top-1 Accuracy (%)
CifarNet	Conventional	217.32	78.2
	Deep Reuse	154.44	73.2 ~ 76.1
	TREC	153.92	76.5
ZfNet	Conventional	3557.32	80.1
	Deep Reuse	814.03	72.5 ~ 76.6
	TREC	814.01	78.9
Vanilla SqueezeNet	Conventional	1639.51	83.5
	Deep Reuse	328.97	79.8 ~ 81.9
	TREC	327.90	83.0
SqueezeNet + Complex Bypass	Conventional	1998.86	85.3
	Deep Reuse	543.71	80.5 ~ 83.1
	TREC	544.03	84.6
ResNet-34 (ImageNet-64×64)	Conventional	4242.77	52.6
	Deep Reuse	1379.26	46.7 ~ 49.9
	TREC	1378.75	52.2

Thanks for listening

guanjw@ruc.edu.cn