

# Fractional-Overlap Declustered Parity: Evaluating Reliability for Storage Systems

**Huan Ke,**  
Haryadi S. Gunawi,



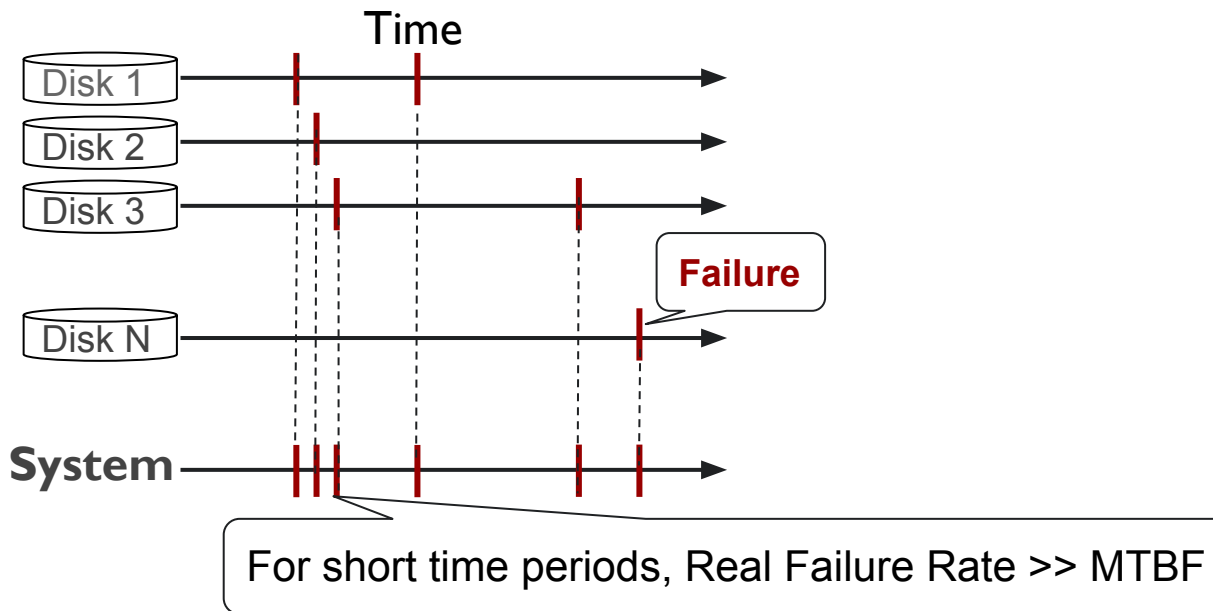
THE UNIVERSITY OF  
**CHICAGO**

Dominic Manno, David Bonnie,  
Bradley W. Settlemyer



# Correlated Failures

Correlated failures within compressed time windows make storage systems highly vulnerable to data loss.



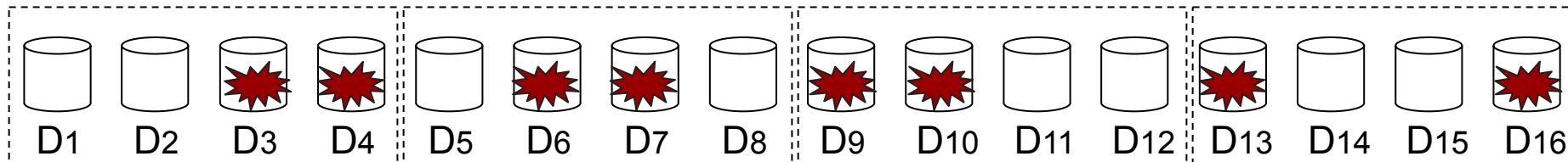
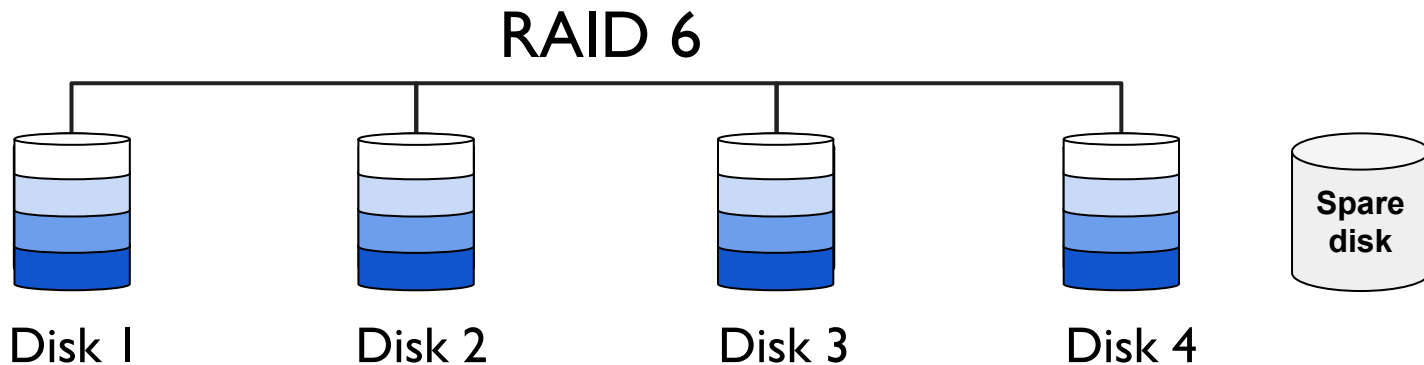
# Failure Models

How do we model correlated failures ...

Types	Models
Poisson Failures	$Poisson(\frac{1}{MTBF})$
Exponential Failures	$Exp(\frac{1}{MTBF})$
Batch Failures	$Exp(\frac{1}{MTBF}) \& Exp(\frac{1}{0.1MTBF})$

# Traditional RAID

RAID (Redundant Array of Inexpensive Disks)



# Declustered Parity (DP)

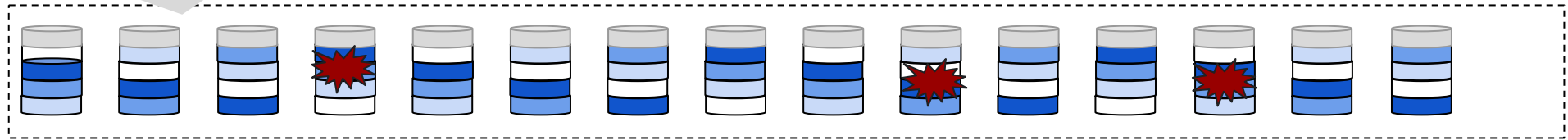
Data/parity are declustered or spread across all disks.

distributed spare space

parallel reads/writes

ZFS dRAID

GridRAID



Spare  
disk

**The probability of data loss is 100%**

# Motivations



Fault Tolerance

- Slower reconstruction

Traditional RAID

Declustered Parity

Rebuild Performance



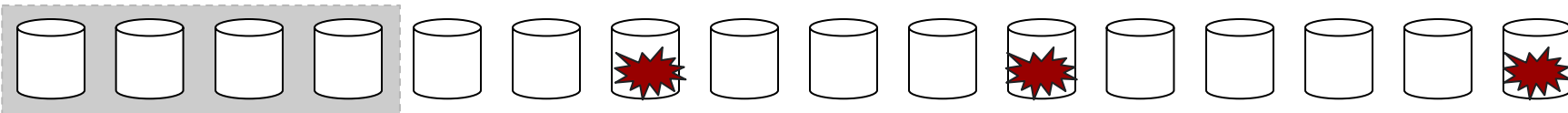
- Lower fault tolerance


How the interactions between **fault tolerance** and **rebuild performance** together impact system reliability is still unclear.


# Fractional Overlap Declustered Parity


FODP, a flexible tool to explore the middle space between fault tolerance and rebuild performance.


D1 D2 D3 D4 D5 D6 D7 D8 D9 D10 D11 D12 D13 D14 D15 D16

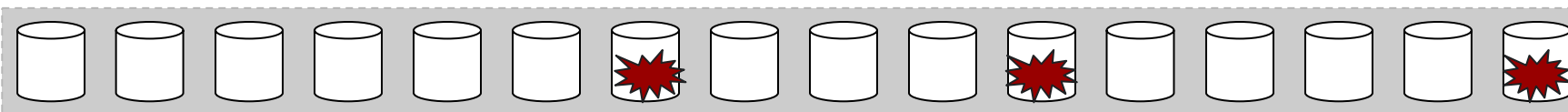


 Flexible rebuild performance

 Uniform data distribution

 Adjustable failure domains

 Higher fault tolerance

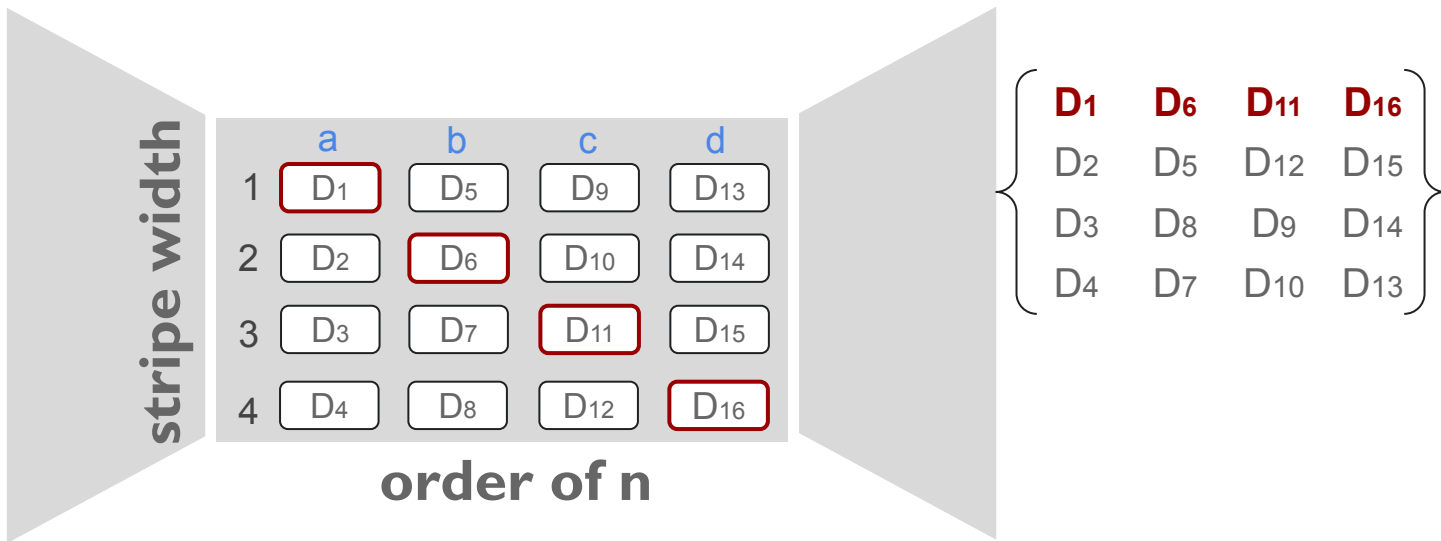


# FODP Construction

## Latin square of order n:

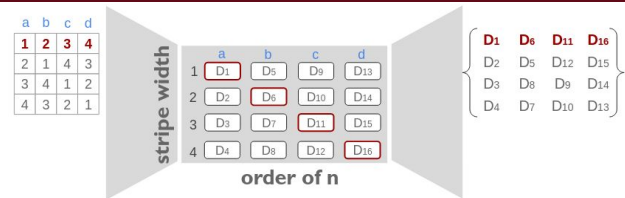
- a  $n \times n$  array over  $n$  elements and each element appears once in each row and column.

	a	b	c	d
1	1	2	3	4
2	2	1	4	3
3	3	4	1	2
4	4	3	2	1





# Overlap fraction



Each latin square corresponds to  $n$  disk subsets that cover the whole disk matrix.

- Each disk has  $(\text{stripe-width} - 1)$  **overlaps** within a disk subset.

Overlap fraction  $\lambda$  for each disk:

$$\lambda = \frac{\text{overlaps}}{N-1}$$

	RAID	FODP	SODP	DP
	$\lambda_{min}$	$\lambda < 1$	$\lambda = 1$	$\lambda > 1$
Rebuild Perf	L	M	H	H
Fault Tolerance	H	H	M	L

# Mutually Orthogonal Latin Squares

Two latin squares are **mutually orthogonal**:

- Any order pair of entries from each latin square in the same row and column occurs exactly once.

1	2	3	4
2	1	4	3
3	4	1	2
4	3	2	1

 + 

1	3	4	2
2	4	3	1
3	1	2	4
4	2	1	3

 → 

1,1	2,3	3,4	4,2
2,2	1,4	4,3	3,1
3,3	4,1	1,2	2,4
4,4	3,2	2,1	1,3

- With any given order of  $n$ , there can be at most  $(n-1)$  mutually orthogonal latin squares (MOLS).

# MOLS in FODP

$L_1$

	a	b	c	d
1	1	2	3	4
2	2	1	4	3
3	3	4	1	2
4	4	3	2	1

$L_2$

1	1	3	4	2
2	2	4	3	1
3	3	1	2	4
4	4	2	1	3

$L_3$

1	1	4	2	3
2	2	3	1	4
3	3	2	4	1
4	4	1	3	2

	a	b	c	d
1	D <sub>1</sub>	D <sub>5</sub>	D <sub>9</sub>	D <sub>13</sub>
2	D <sub>2</sub>	D <sub>6</sub>	D <sub>10</sub>	D <sub>14</sub>
3	D <sub>3</sub>	D <sub>7</sub>	D <sub>11</sub>	D <sub>15</sub>
4	D <sub>4</sub>	D <sub>8</sub>	D <sub>12</sub>	D <sub>16</sub>

D <sub>1</sub>	D <sub>6</sub>	D <sub>11</sub>	D <sub>16</sub>
D <sub>2</sub>	D <sub>5</sub>	D <sub>12</sub>	D <sub>15</sub>
D <sub>3</sub>	D <sub>8</sub>	D <sub>9</sub>	D <sub>14</sub>
D <sub>4</sub>	D <sub>7</sub>	D <sub>10</sub>	D <sub>13</sub>

$$\lambda = \frac{3}{15}$$

D <sub>1</sub>	D <sub>7</sub>	D <sub>12</sub>	D <sub>14</sub>
D <sub>2</sub>	D <sub>8</sub>	D <sub>11</sub>	D <sub>13</sub>
D <sub>3</sub>	D <sub>5</sub>	D <sub>10</sub>	D <sub>16</sub>
D <sub>4</sub>	D <sub>6</sub>	D <sub>9</sub>	D <sub>15</sub>

$$\lambda = \frac{6}{15}$$

D <sub>1</sub>	D <sub>8</sub>	D <sub>10</sub>	D <sub>15</sub>
D <sub>2</sub>	D <sub>7</sub>	D <sub>9</sub>	D <sub>16</sub>
D <sub>3</sub>	D <sub>6</sub>	D <sub>12</sub>	D <sub>13</sub>
D <sub>4</sub>	D <sub>5</sub>	D <sub>11</sub>	D <sub>14</sub>

$$\lambda = \frac{9}{15}$$

# Trade-offs in FODP

FODP gives us the flexibility to explore the trade-offs between fault tolerance and rebuild performance.

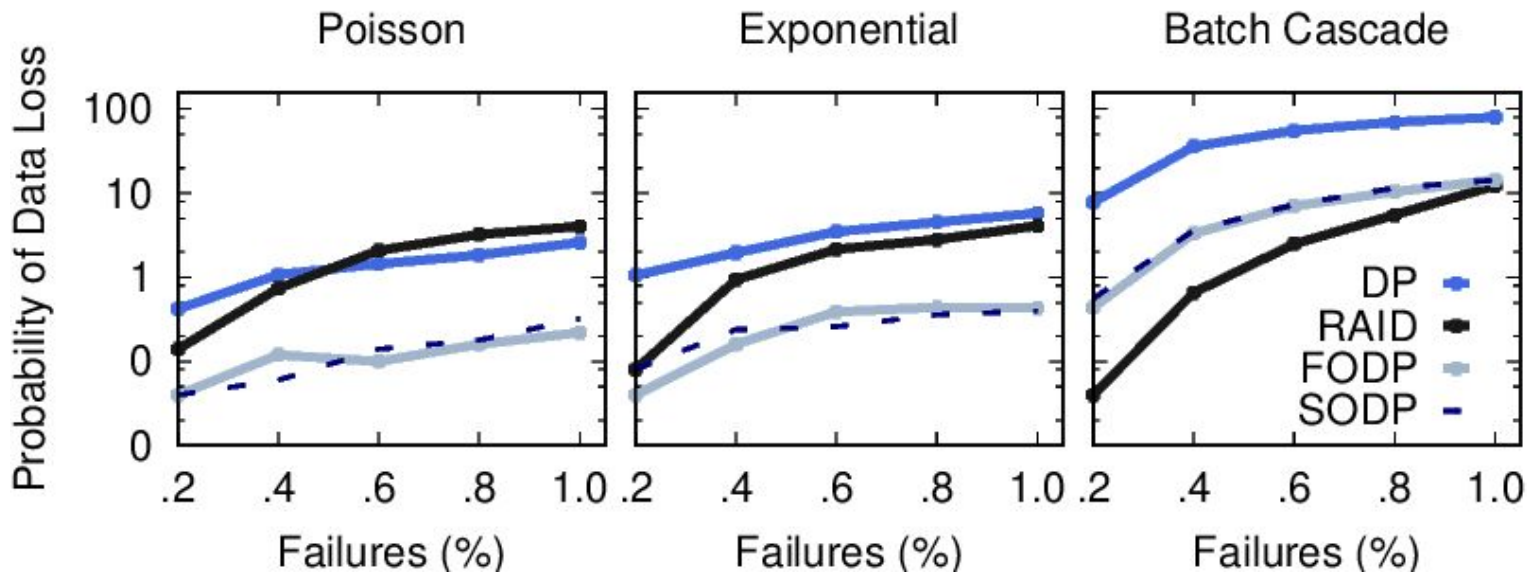
- ❑ The lower  $\lambda$  is, the more failures that can be tolerated.
- ❑ The larger  $\lambda$  is, the more overlaps can be used for rebuilds.

**FODP+1**

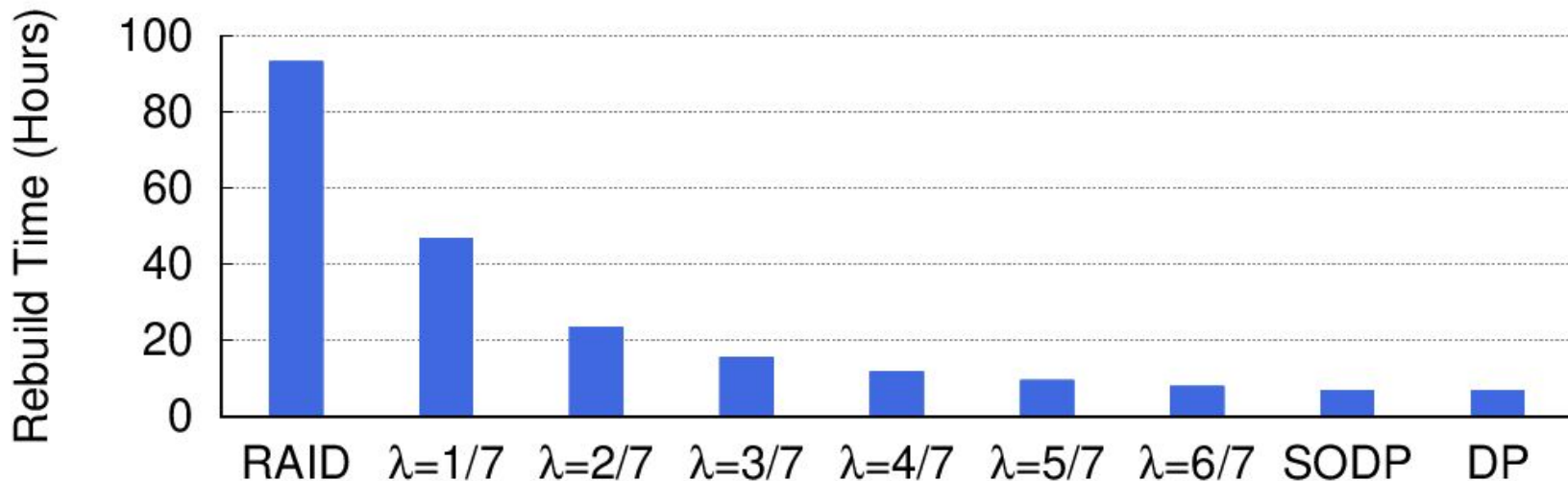
**If data loss occurs, FODP loses more data than DP**

# Impact of Failures

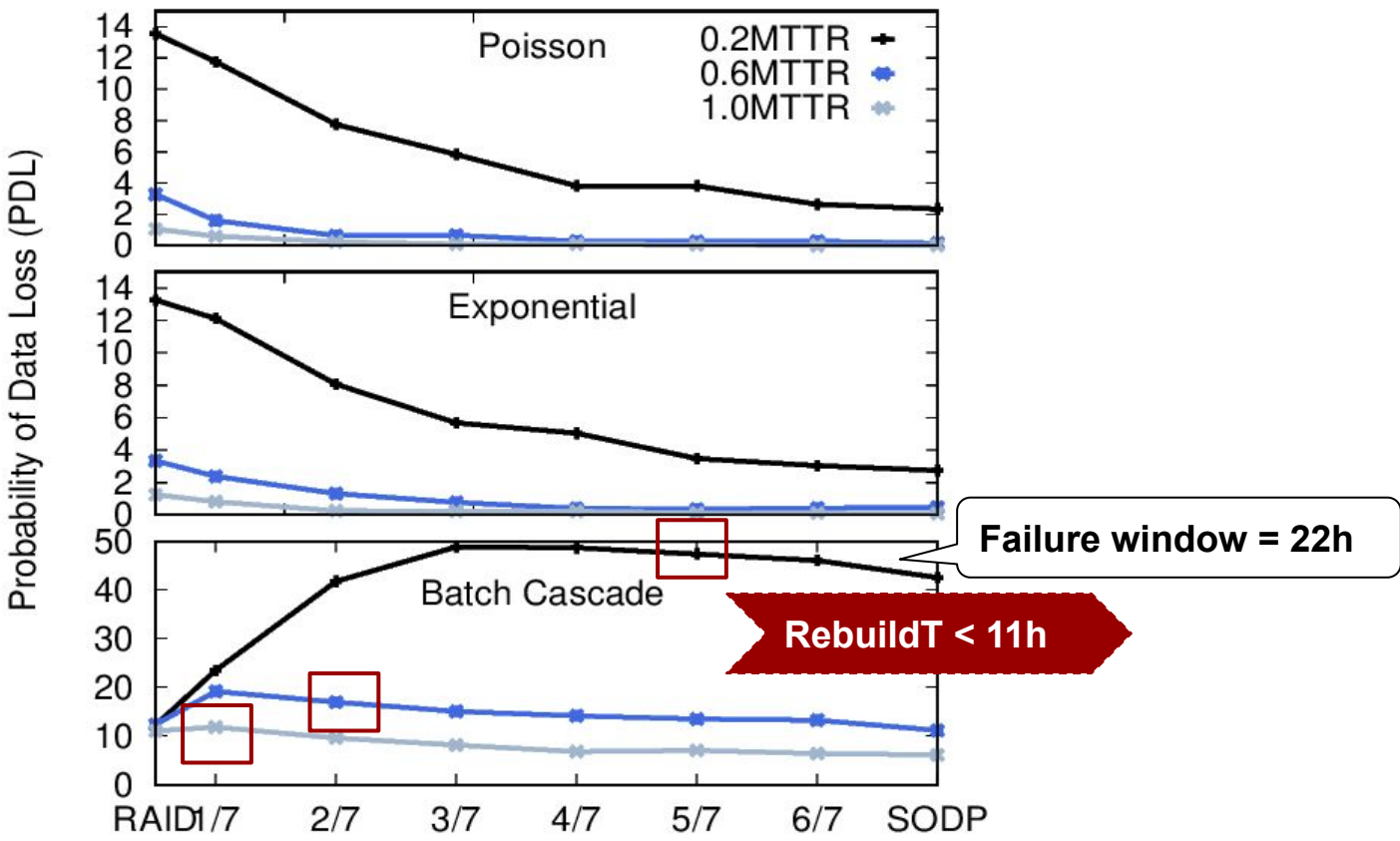
Assume  $MTBF = 0.5 MTTR$  in Campaign system with 11+2 configurations within each server.



# Impact of Overlap Fraction



# Impact of Overlap Fraction



# FODP Conclusion

*“Why should we address correlated failures?”*

Storage systems are becoming **larger** and **denser** and failures are increasingly **correlated in time!**

**FODP**, a flexible tool to study and explore rebuild performance and failure domains in systems.

**FODP-Plus-One**, reducing the magnitude of data loss by adding a layer of parity on top of FODP stripes.





# Thank you!

# Questions?



<http://ucare.cs.uchicago.edu>