# Regenerating Codes for Errers and Erasures in Distributed Storage

Nihar Shah
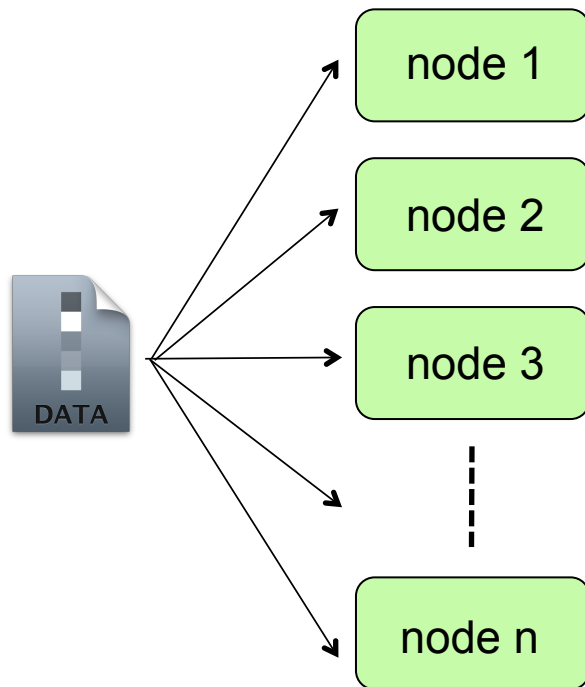
joint work with

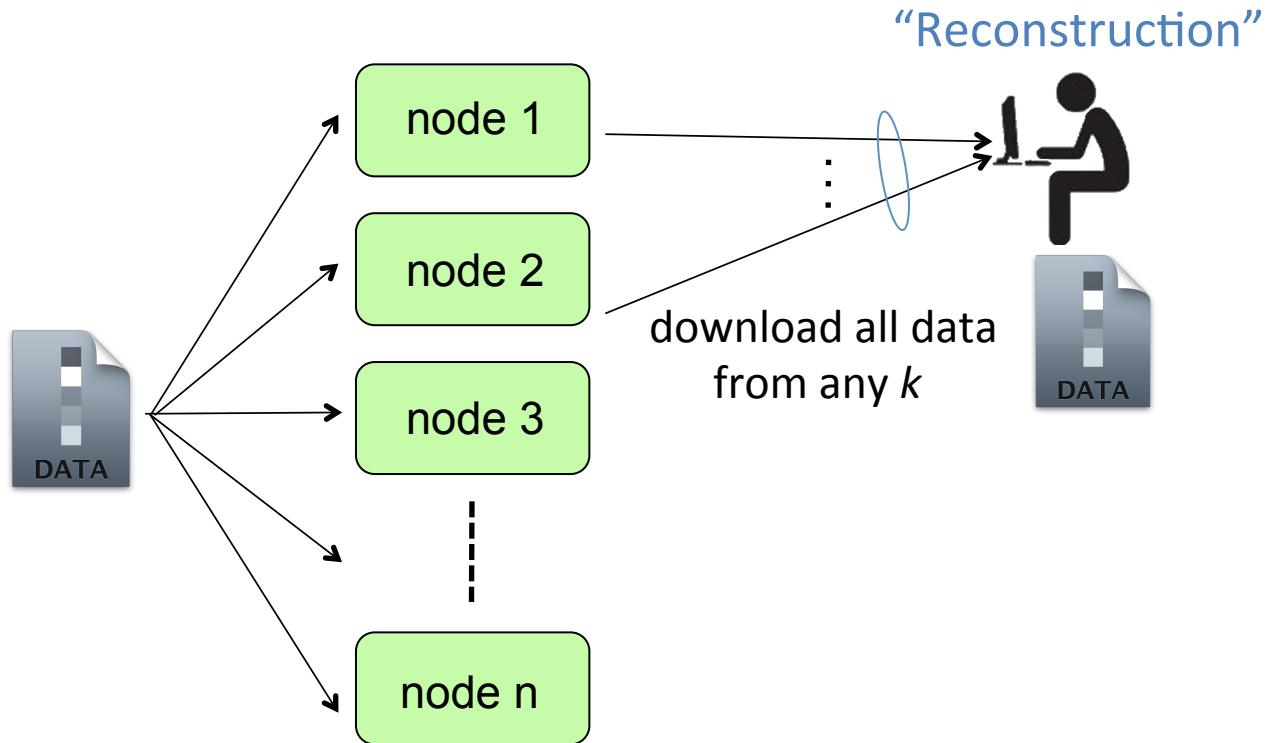K V Rashmi ,  Kannan Ramchandran ,  P Vijay Kumar

# Regenerating Codes

- Provide reliability
- Efficient repair



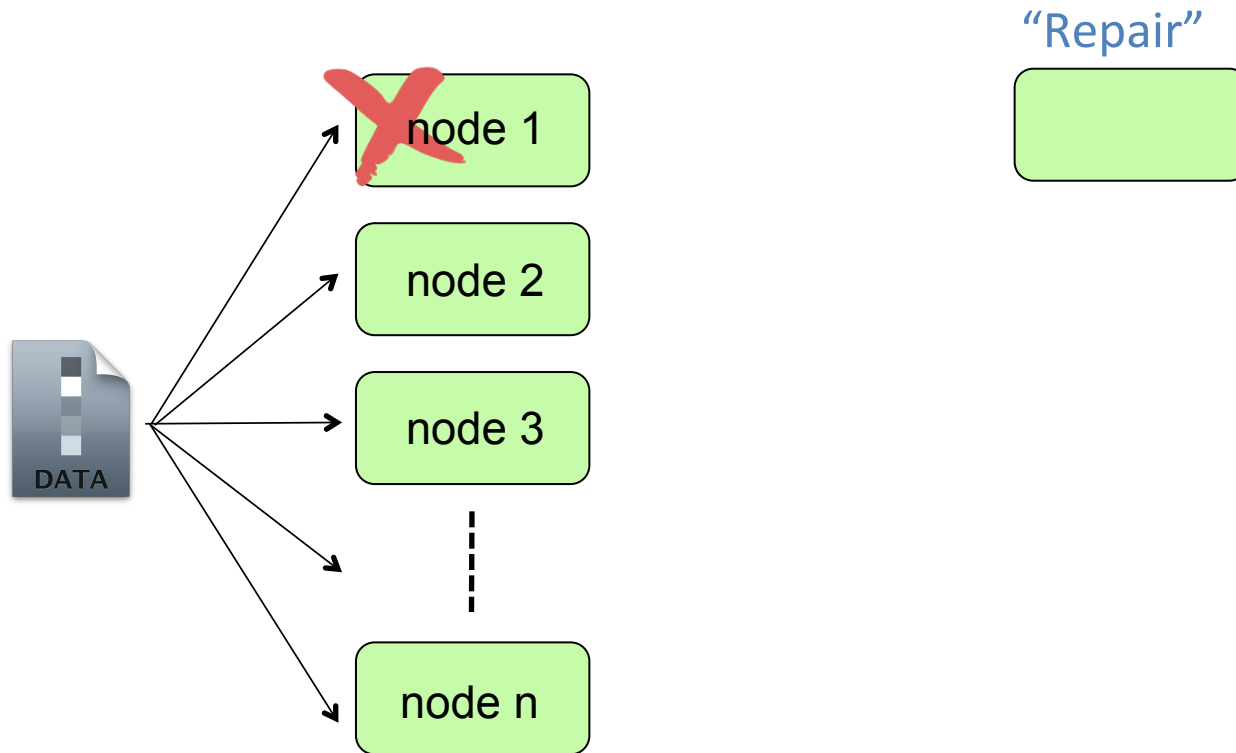Introduced by Dimakis et al. (Infocom 07, IT-Transactions 10)

# Regenerating Codes

- Provide reliability: recover data from any k nodes
- Efficient repair



"Reconstruction"

node 1

node 2

node 3

node n

DATA
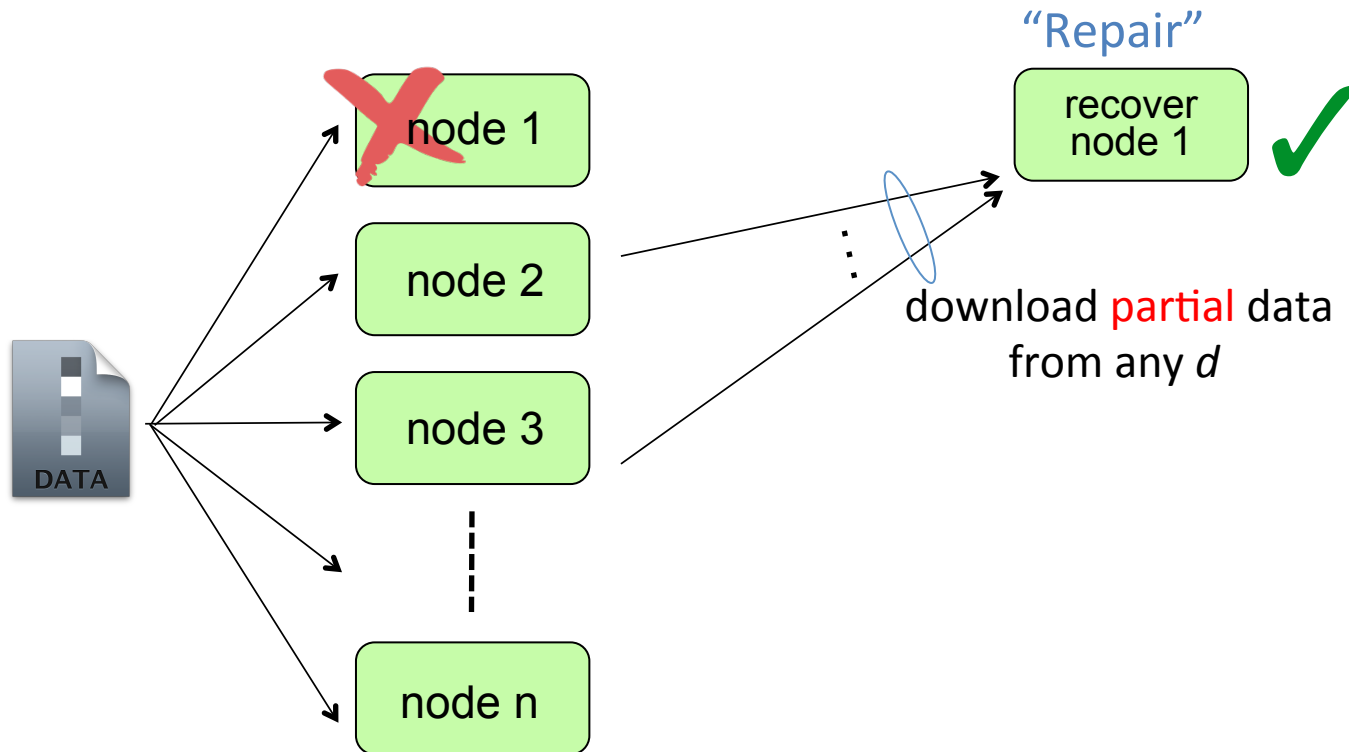
download all data from any *k*
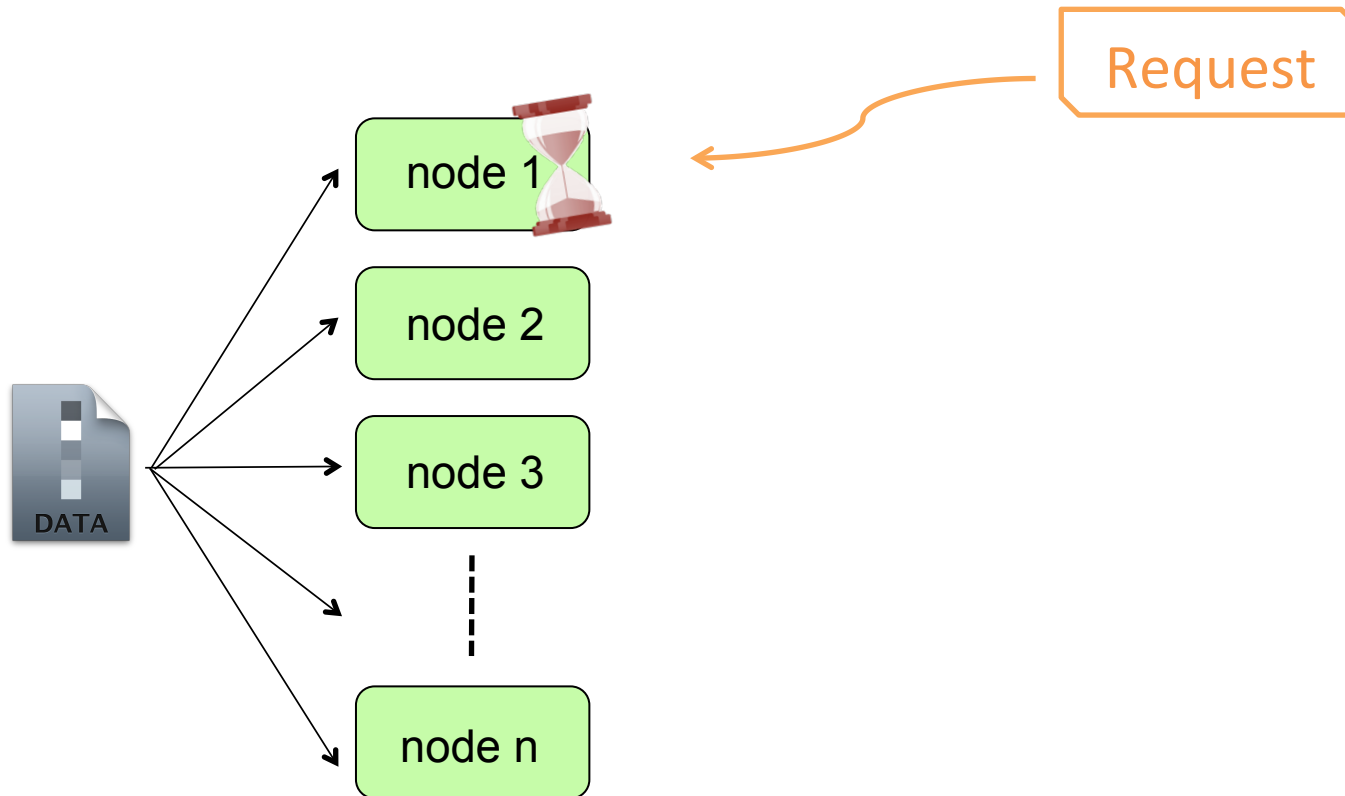
DATA

# Regenerating Codes

- Provide reliability: recover data from any k nodes
- Efficient repair: small network bandwidth

# Regenerating Codes

- Provide reliability: recover data from any k nodes
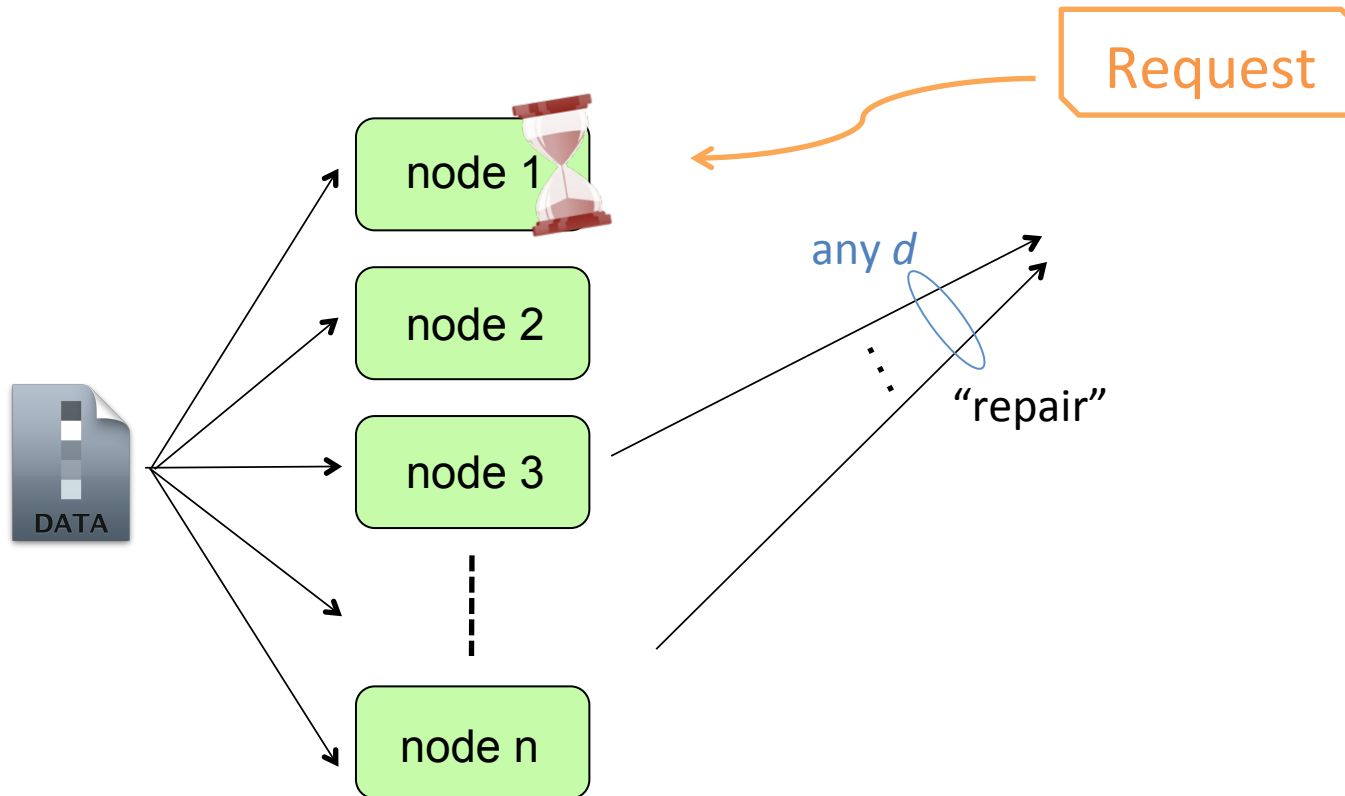- Efficient repair: small network bandwidth

# Regenerating Codes

- Provide reliability: recover data from any k nodes

- Efficient repair: small network bandwidth
  - (Fast) degraded reads

# Regenerating Codes

- Provide reliability: recover data from any k nodes
- Efficient repair: small network bandwidth
  - (Fast) degraded reads

# Explicit Regenerating Code Constructions

Minimum Bandwidth Point (MBR):
    Rashmi et al.'09, Rashmi et al.'11

Minimum Storage Point (MSR):
    Rashmi et al.'09, Shah et al. '09,
    Suh et al.'10, Rashmi et al.'11, Cadambe et al.'11,
    Tamo et al.'11, Wang et al.'11, Papailiopoulos et al.'11
….
Cooperative repair, adaptive repair, flexible repair…

# Explicit Regenerating Code Constructions

Minimum Bandwidth Point (MBR):
     Rashmi et al.'09, Rashmi et al.'11
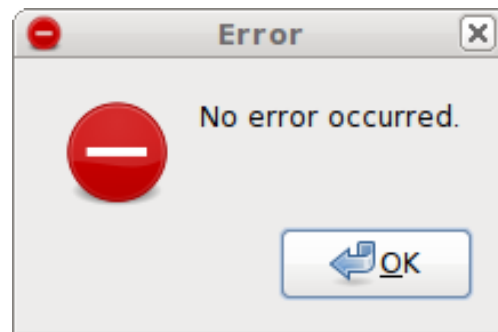
Minimum Storage Point (MSR):
     Rashmi et al.'09, Shah et al. '09,
     Suh et al.'10, Rashmi et al.'11, Cadambe et al.'11,
     Tamo et al.'11, Wang et al.'11, Papailiopoulos et al.'11

Cooperative repair, adaptive repair, flexible repair…

Assume an error/erasure-free setting



In this talk: errors and erasures…

# Motivation I: FEC for Network Errors

# Motivation II: Improve Latency

# Motivation II: Improve Latency

Request

node 1

node 2

node 3

node n

$d+2$

Download from first d who respond
(treat remaining two as erasures)

# Motivation III: Security

node 1

node 2

node 3

node n

- Compromised nodes transmit erroneous data

# Motivation III: Security



- Compromised nodes transmit erroneous data
- Errors may propagate during repair operations

# Motivation III: Security



- Compromised nodes transmit erroneous data
- Errors may propagate during repair operations

# Motivation III: Security



- Compromised nodes transmit erroneous data
- Errors may propagate during repair operations
- Outer Bounds: Pawar et al. '11

# Motivation IV: Theoretical Interest

- Establish capacity of such systems

# Code Constructions

Explicit codes performing error and erasure correction for

- Minimum Bandwidth (MBR):    all parameters
- Minimum Storage (MSR):     all $[n, k, d \geq 2k-2]$

These codes achieve (and establish) capacity.

# Recipe



Encode using (error-free)
Product-Matrix code[1]

[1] Rashmi, Shah, Kumar, IEEE Transactions on Information Theory, 2011

# Recipe



Encode using (error-free)
Product-Matrix code[1]

$d + \Delta$

download from
additional nodes

[1] Rashmi, Shah, Kumar, IEEE Transactions on Information Theory, 2011

# Toy Example: Minimum Bandwidth Code

- Tolerate any 3 node failures

node 1

node 2

node 3

node 4

node 5

# Toy Example: Minimum Bandwidth Code

- Data = {a, b, c}

node 1

node 2

node 3

node 4

node 5

# Toy Example: Minimum Bandwidth Code

- Data = {a, b, c}

$$\begin{bmatrix} a & b \\ b & c \end{bmatrix}$$

Message
matrix

node 1

node 2

node 3

node 4

node 5

# Toy Example: Minimum Bandwidth Code

- Data = {a, b, c}

[ 1   0 ]

[ 0   1 ]

[ 1   1 ]     $\begin{bmatrix} a & b \\ b & c \end{bmatrix}$

[ 1   2 ]     Message matrix

[ 1   3 ]

Encoding vectors

node 1

node 2

node 3

node 4

node 5

# Toy Example: Minimum Bandwidth Code

- Data = {a, b, c}

[ 1    0 ]

[ 0    1 ]

[ 1    1 ] x $\begin{bmatrix} a & b \\ b & c \end{bmatrix}$

Message matrix

[ 1    2 ]

[ 1    3 ]

Encoding vectors

| a | b | node 1 |
| b | c | node 2 |
| a + b | b + c | node 3 |
| a + 2b | b + 2c | node 4 |
| a + 3b | b + 3c | node 5 |

# Toy Example: Minimum Bandwidth Code

- Can recover data from any 2 nodes (in absence of errors/erasures)

$[\ 1\ \ \ 0\ ]$

$[\ 0\ \ \ 1\ ]$

$[\ 1\ \ \ 1\ ]\ \times\ \begin{bmatrix} a & b \\ b & c \end{bmatrix}$

Message matrix

$[\ 1\ \ \ 2\ ]$

$[\ 1\ \ \ 3\ ]$

Encoding vectors

| a | b | node 1 |
|---|---|---|
| b | c | node 2 |
| a + b | b + c | node 3 |
| a + 2b | b + 2c | node 4 |
| a + 3b | b + 3c | node 5 |

a , b , c

# Toy Example: Minimum Bandwidth Code

- Optimal Repair (in absence of errors/erasures)

# Toy Example: Minimum Bandwidth Code

- Optimal Repair (in absence of errors/erasures)

| node 1 | a ✗ b |
| node 2 | b | c |
| node 3 | a + b | b + c |
| node 4 | a + 2b | b + 2c |
| node 5 | a + 3b | b + 3c |

$$\text{node 2} \times \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$\text{node 3} \times \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

Encoding vector
$[ 1 \quad 0 ]$

# Toy Example: Minimum Bandwidth Code

- Optimal Repair (in absence of errors/erasures)

| node 1 | a ✗ b |
|--------|-------|
| node 2 | b | c |
| node 3 | a + b | b + c |
| node 4 | a + 2b | b + 2c |
| node 5 | a + 3b | b + 3c |

$\text{node 2} \times \begin{bmatrix} 1 \\ 0 \end{bmatrix}$  →  $b$

$\text{node 3} \times \begin{bmatrix} 1 \\ 0 \end{bmatrix}$  →  $a + b$

Encoding vector
$[\ 1 \quad 0\ ]$

# Toy Example: Minimum Bandwidth Code

- Optimal Repair (in absence of errors/erasures)



node 1: a | b (marked with red X)

node 2: b | c $\times \begin{bmatrix} 1 \\ 0 \end{bmatrix}$

node 3: a + b | b + c $\times \begin{bmatrix} 1 \\ 0 \end{bmatrix}$

node 4: a + 2b | b + 2c

node 5: a + 3b | b + 3c

b

a + b

a | b (marked with green check)

Encoding vector
[ 1    0 ]

# Toy Example: Minimum Bandwidth Code

- Optimal Repair (in absence of errors/erasures)



Encoding vector
[ 1    0 ]

# Toy Example: Minimum Bandwidth Code

- Optimal Repair (in absence of errors/erasures)

# Toy Example: Minimum Bandwidth Code

- Optimal Repair (in absence of errors/erasures)

# Toy Example: Minimum Bandwidth Code

- For repair resilient to one erasure: connect to one additional node

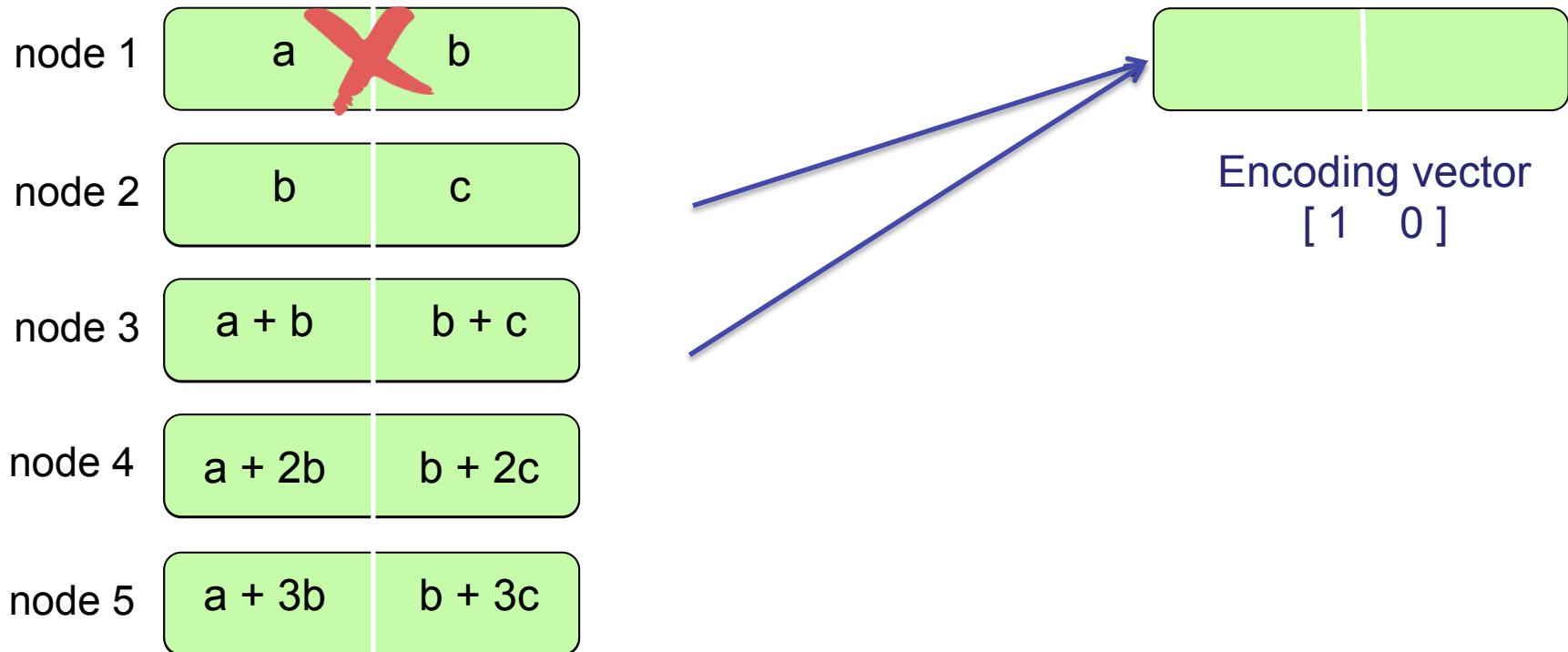| | |
|---|---|
| node 1 | a $\quad$ **X** $\quad$ b |
| node 2 | b $\quad\quad$ c |
| node 3 | a + b $\quad\quad$ b + c |
| node 4 | a + 2b $\quad\quad$ b + 2c |
| node 5 | a + 3b $\quad\quad$ b + 3c |

Encoding vector
[ 1 $\quad$ 0 ]

# Toy Example: Minimum Bandwidth Code



node 1 | a | b

node 2 | b | c | $\times \begin{bmatrix} 1 \\ 0 \end{bmatrix}$

node 3 | a + b | b + c | $\times \begin{bmatrix} 1 \\ 0 \end{bmatrix}$

node 4 | a + 2b | b + 2c

node 5 | a + 3b | b + 3c | $\times \begin{bmatrix} 1 \\ 0 \end{bmatrix}$

Encoding vector
[ 1    0 ]

# Toy Example: Minimum Bandwidth Code

node 1 | a ✗ b

node 2 | b | c | × $\begin{bmatrix} 1 \\ 0 \end{bmatrix}$

node 3 | a + b | b + c | × $\begin{bmatrix} 1 \\ 0 \end{bmatrix}$

node 4 | a + 2b | b + 2c

node 5 | a + 3b | b + 3c | × $\begin{bmatrix} 1 \\ 0 \end{bmatrix}$

b

a + b

a + 3b

a ✓ b

Encoding vector
[ 1    0 ]

[3,2] MDS code in {a, b}
Can correct one arbitrary erasure

(uses special structure
of product-matrix codes)

# Toy Example: Minimum Bandwidth Code

- For repair resilient to one error: connect to two additional nodes



node 1: a | b

node 2: b | c

node 3: a + b | b + c

node 4: a + 2b | b + 2c

node 5: a + 3b | b + 3c

Encoding vector
[ 1    0 ]

# Toy Example: Minimum Bandwidth Code



node 1   a   b

node 2   b   c   × $\begin{bmatrix} 1 \\ 0 \end{bmatrix}$

node 3   a + b   b + c   × $\begin{bmatrix} 1 \\ 0 \end{bmatrix}$

node 4   a + 2b   b + 2c   × $\begin{bmatrix} 1 \\ 0 \end{bmatrix}$

node 5   a + 3b   b + 3c   × $\begin{bmatrix} 1 \\ 0 \end{bmatrix}$

b

a + b

a + 2b + e

a + 3b

a   b

Encoding vector
[ 1   0 ]

[4,2] MDS code in {a, b}
Can correct one arbitrary error

# Toy Example: Minimum Bandwidth Code

- Data-reconstruction in presence of errors/erasures



| | | |
|---|---|---|
| node 1 | a | b |
| node 2 | b | c |
| node 3 | a + b | b + c |
| node 4 | a + 2b | b + 2c |
| node 5 | a + 3b | b + 3c |

a , b , c

"MDS" property corrects errors/erasures

(any 'k' property)

# General Algorithm

- Encode and store using Product-Matrix code for the error/erasure-free case

- For resilience from $s$ errors and $t$ erasures, connect to
  - $d+2s+t$ for repair
  - $k+2s+t$ for reconstruction

- Helping nodes oblivious to resilience requirements

- Resulting data passed is MDS



$g_1^T M$

$g_2^T M$

$g_3^T M$

$g_n^T M$

$g_1^T M h_f$

$g_2^T M h_f$

$g_{d+2b}^T M h_f$

$M h_f \rightarrow g_f^T M$

[d+2s+t, d] MDS code
Corrects s errors, t erasures

# Analysis I: Secure-capacity

- Meets outer bound (Pawar et al. '11)

$$B \leq \sum_{i=2s+t}^{k-1} \min(\alpha, (d-i)\beta)$$

Establishes the capacity of these systems

- Related to node compromise in network coding
  - not very well understood in general
  - here: practical, explicit algorithm achieving capacity

# Analysis II: Encoding Independent of Resiliency Requirements

- Meets outer bound

$$B \leq \sum_{i=2s+t}^{k-1} \min(\alpha, (d-i)\beta)$$

simultaneously for all s (#errors) and t (#erasures)

"Universal Resilience"

# Analysis II: Encoding Independent of Resiliency Requirements

- Meets outer bound

$$B \leq \sum_{i=2s+t}^{k-1} \min(\alpha, (d-i)\beta)$$

  simultaneously for all s (#errors) and t (#erasures)

  ┌─────────────────────────┐
  │   "Universal Resilience"  │
  └─────────────────────────┘

- Can choose different protection level for each instance of repair or reconstruction
  - handle time-varying channels/requirements
  - while always remaining optimal !

- Need not design for worst case
  - saves resources

# Analysis III: What about other regenerating codes ?

- Our codes exploit structure of the Product-Matrix framework

Why Product-Matrix framework ?

- Higher connectivity available (d<n-1)
- Data passed by a node independent of other helpers

# Analysis III: What about other regenerating codes ?

Necessary and sufficient:

- Higher connectivity available (d<n-1)
- Data passed by a node independent of other helpers


- Other existing codes[1] restrict number of nodes to n=d+1
  - thus, not applicable in this setting

[1] The high-rate MSR (d=k+1) explicit codes by Rashmi et al. (Allerton '09) do not impose this restriction, however, they perform only approximately-exact repair

# Summary

- Explicit codes for correcting errors and erasures
  - Employing product-matrix framework

- Achieve and establish capacity
  - Related to network-coding with compromised nodes

- Universal resilience
  - Encoding independent of error protection requirements

- Necessary & sufficient conditions for any regenerating code

- Open: capacity in presence of errors/erasures for
  - MSR when redundancy $< \left(2 - \frac{1}{k}\right)$
  - Interior points

# Thanks!

# Backup Slides

# Product-Matrix Codes

- ## Completely solves
  - MBR for all parameters
  - MSR for redundancy $\geq \left(2 - \frac{1}{k}\right)$

- ## Scalable
  - n  independent of all other parameters
  - Only construction supporting arbitrary #nodes

- ## Decentralized
  - Can connect to *any* subset of 'd' nodes for repair

- ## Ready to go
  - Can use (already existing) Reed-Solomon encoders/decoders for implementation

- ## Optimal

# Product-Matrix Framework

Explicit MBR codes for all $n$, $k$, $d$
Explicit MSR codes for all $n$, $k$, $d \in [2k - 2, n - 1]$

$$\underbrace{C}_{n \times \alpha} = \underbrace{\Psi}_{n \times d} \underbrace{M}_{d \times \alpha}$$

- $C$ : Code matrix
  - Every row represents one node
  - $\alpha$ symbols stored in $i^{th}$ node are $\underline{\psi}_i^t M$

- $\Psi$ : Encoding matrix
  - Fixed apriori

- $M$ : Message matrix
  - Contains the $B$ source symbols, with some symbols possibly repeated

# Product-matrix MBR code

- $\alpha = d\beta$ and $B = \frac{k(k+1)}{2} + k(d-k)$

- Code $\underbrace{C}_{n \times \alpha} = \underbrace{\Psi}_{n \times d}\,\underbrace{M}_{d \times \alpha}$

- Message matrix $\underbrace{M}_{d \times d} = \begin{bmatrix} \underbrace{S}_{k \times k} & \underbrace{T}_{k \times (d-k)} \\[1em] \underbrace{T^t}_{(d-k) \times k} & \underbrace{0}_{k \times (d-k)} \end{bmatrix}$

  $S$ symmetric ($\Rightarrow M$ symmetric)

- Encoding matrix $\underbrace{\Psi}_{n \times d} = \begin{bmatrix} \underbrace{\Phi}_{n \times k} & \underbrace{\Delta}_{n \times (d-k)} \end{bmatrix}$

  $\Phi$: any $k$ rows linearly independent
  $\Psi$: any $d$ rows linearly independent
  e.g., $\Psi$ is a Vandermonde or Cauchy matrix

# Product-matrix MBR code : Node repair

Replacement node $f$ needs: $\underline{\psi}_f^t M$

Helper node $i$ stores: $\underline{\psi}_i^t M$

---

Helper node $i$ passes: $\underline{\psi}_i^t M \underline{\psi}_f$

From $d$ nodes $\quad \downarrow$

$\Psi_{rep} M \underline{\psi}_f$

($\Psi_{rep}$ is $d \times d$, invertible)

$\downarrow$

$M \underline{\psi}_f$

($M$ is symmetric)

$\downarrow$

$\underline{\psi}_f^t M$

$$\underbrace{C}_{n \times \alpha} = \underbrace{\Psi}_{n \times d} \underbrace{M}_{d \times alpha}$$

$$\underbrace{M}_{d \times d} = \begin{bmatrix} \underbrace{S}_{k \times k} & \underbrace{T}_{k \times (d-k)} \\ \underbrace{T^t}_{(d-k) \times k} & \underbrace{0}_{k \times (d-k)} \end{bmatrix}$$

$S$ symmetric $\Rightarrow M$ symmetric

$$\underbrace{\Psi}_{n \times d} = \begin{bmatrix} \underbrace{\Phi}_{n \times k} & \underbrace{\Delta}_{n \times (d-k)} \end{bmatrix}$$

$\Phi$: any $k$ rows LI
$\Psi$: any $d$ rows LI

# Product-matrix MBR code : Data Reconstruction

Node $i$ passes: $\underline{\psi}_i^t M$

From $k$ nodes $\downarrow$

$\Psi_{\mathrm{DC}} M$

$(\Psi_{\mathrm{DC}} = [\Phi_{\mathrm{DC}} \ \ \Delta_{\mathrm{DC}}]$ is $(k \times d))$

$\downarrow$

$\left[ \ \ \Phi_{\mathrm{DC}} S + \Delta_{\mathrm{DC}} T^t \ \ \ \ \Phi_{\mathrm{DC}} T \ \ \right]$

$\downarrow$

$\Phi_{\mathrm{DC}}$ is $k \times k$, invertible
Decode $T$

$\downarrow$

Subtract $\Delta_{\mathrm{DC}} T^t$, Decode $S$

---

$\underbrace{C}_{n \times \alpha} = \underbrace{\Psi}_{n \times d} \ \underbrace{M}_{d \times alpha}$

$\underbrace{M}_{d \times d} = \begin{bmatrix} \underbrace{S}_{k \times k} & \underbrace{T}_{k \times (d-k)} \\ \underbrace{T^t}_{(d-k) \times k} & \underbrace{0}_{k \times (d-k)} \end{bmatrix}$

$S$ symmetric $\Rightarrow M$ symmetric

$\underbrace{\Psi}_{n \times d} = \begin{bmatrix} \underbrace{\Phi}_{n \times k} & \underbrace{\Delta}_{n \times (d-k)} \end{bmatrix}$

$\Phi$: any $k$ rows LI
$\Psi$: any $d$ rows LI

# Product-Matrix MSR Codes

- $\alpha = \frac{B}{k}$
- Hence, necessarily MDS (over alphabet $\mathbb{F}_q^\alpha$)
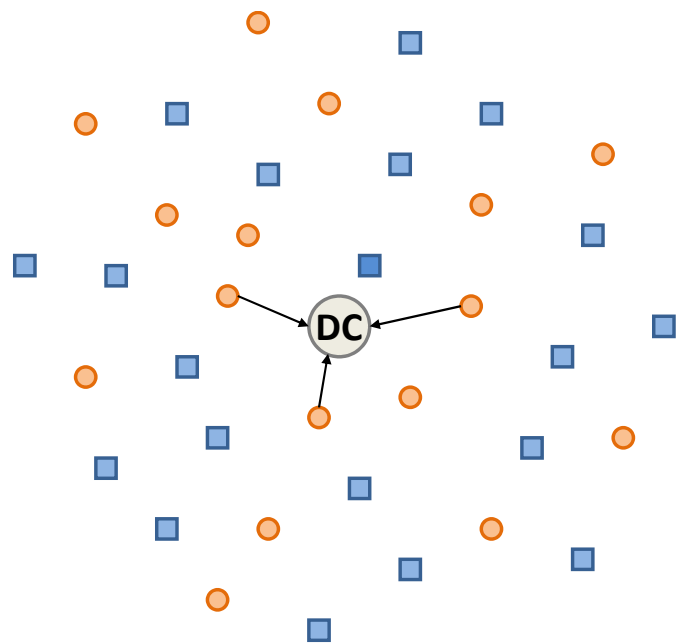- $d = (k-1) + \alpha$     (striping: $\beta = 1$)

- We design codes for all $n$, $k$, $d \in [2k-2, n-1]$
- $C = \Psi M$, $\Psi = [\Phi\ \Lambda\Phi]$, $M = \begin{bmatrix} S_1 \\ S_2 \end{bmatrix}$

- Involves solving multiple simultaneous interference-alignment constraints

# Twin Codes

- New framework to allow use of *any erasure code* and still have efficient repair

- Properties of these constituent codes used during data reconstruction/repair
  - low complexity decoding
  - error detection/correction
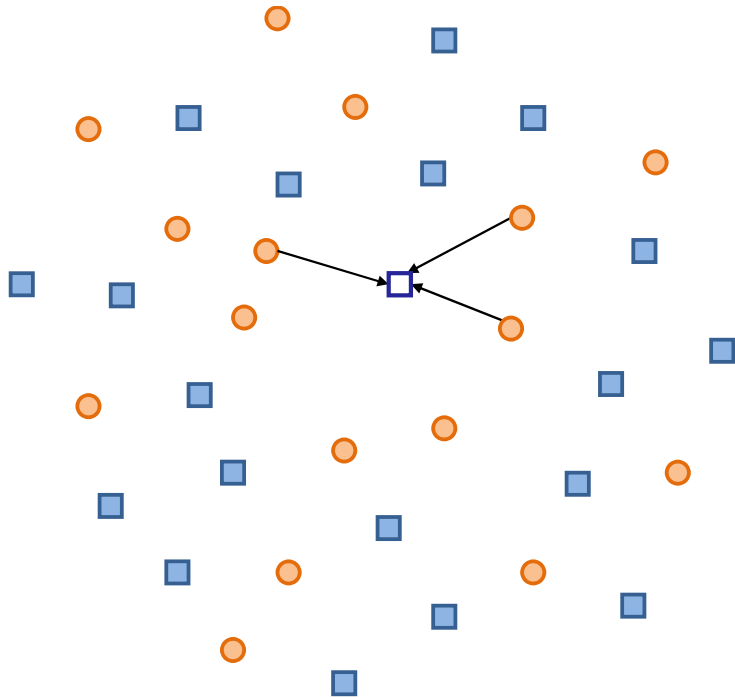  - ratelessness
  - etc.

# Twin Codes



Two types of nodes

- encoded using two different codes

Data reconstruction by connecting to nodes of the *same* type

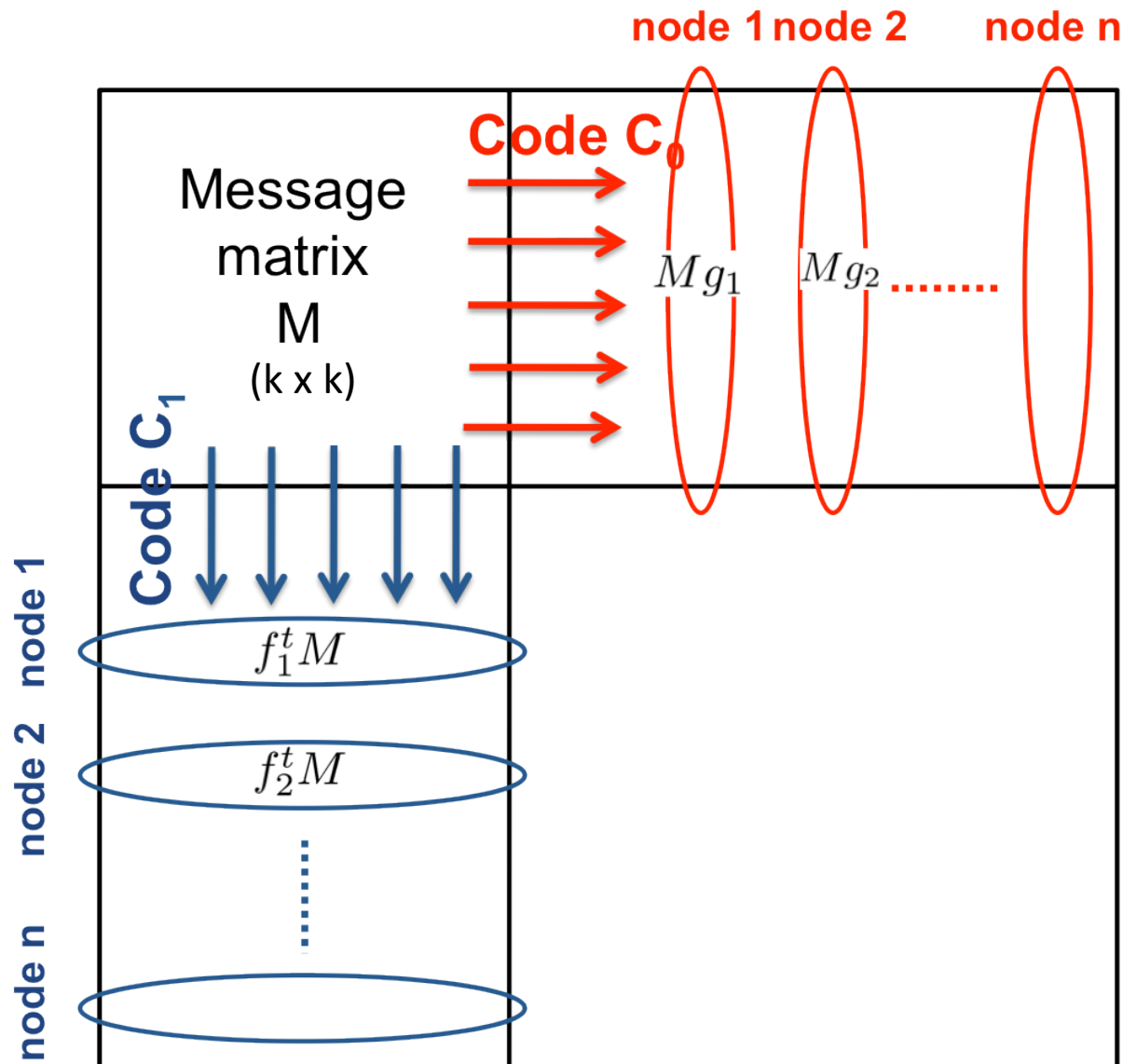**DC** Data collector
■ Type 0 storage node
● Type 1 storage node

# Twin Codes



Node repair by connecting to nodes of the *other* type

□ **Replacement node of Type 0**

# Twin Codes Construction



node 1 node 2    node n

**Code $C_0$**

Message matrix M (k x k)

**Code $C_1$**

$Mg_1$    $Mg_2$ ........

$f_1^t M$

$f_2^t M$

node 1    node 2    node n

$f_i$

# Twin Codes Construction

node 1 node 2        node n

**Code C₀**

Message matrix M (k x k)

$Mg_1$     $Mg_2$ ........
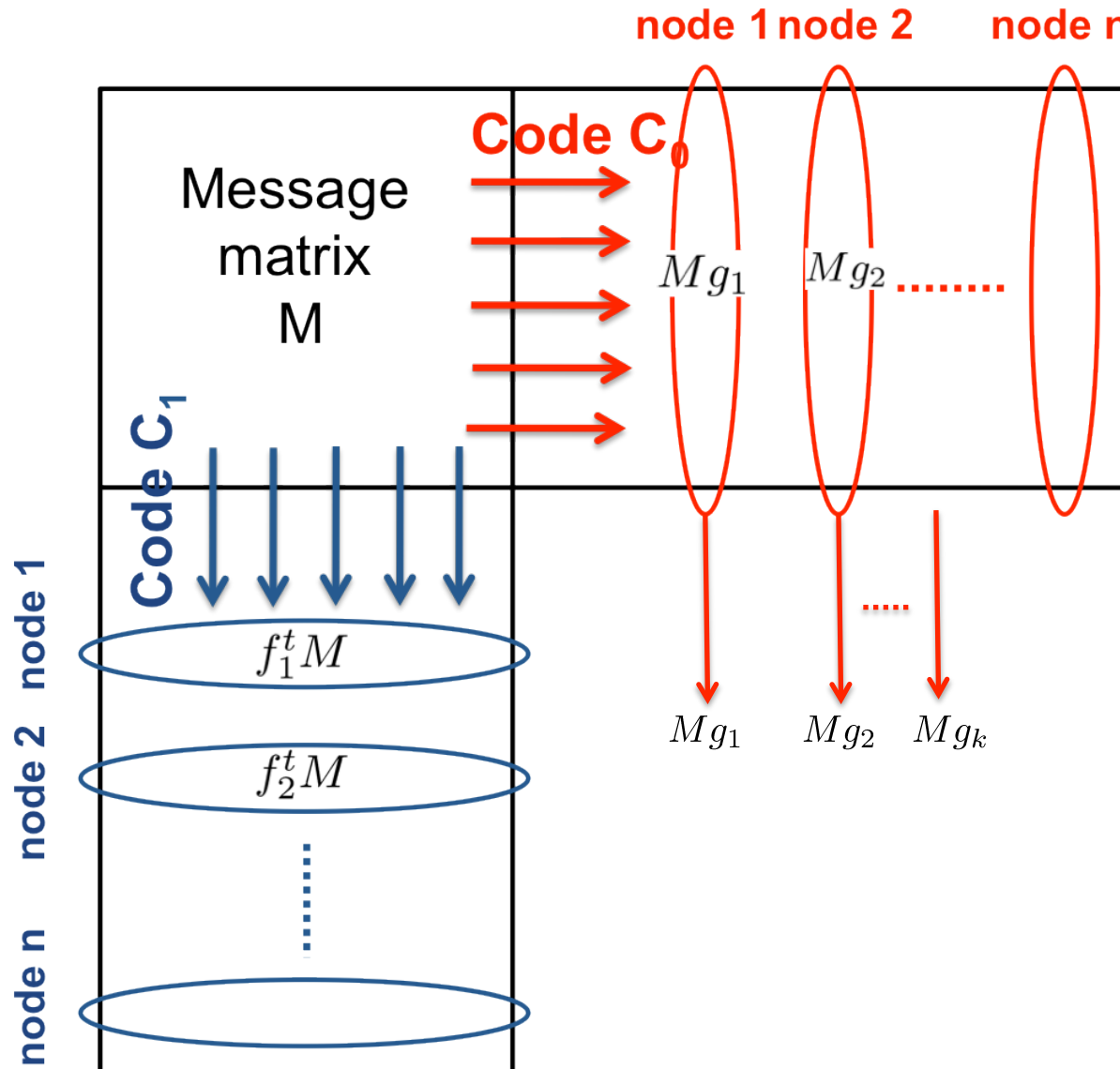
**Code C₁**

$f_1^t M$

$f_2^t M$

node n   node 2   node 1

- any k of the $g_i$ are linearly independent

- any k of the $f_i$ are linearly independent
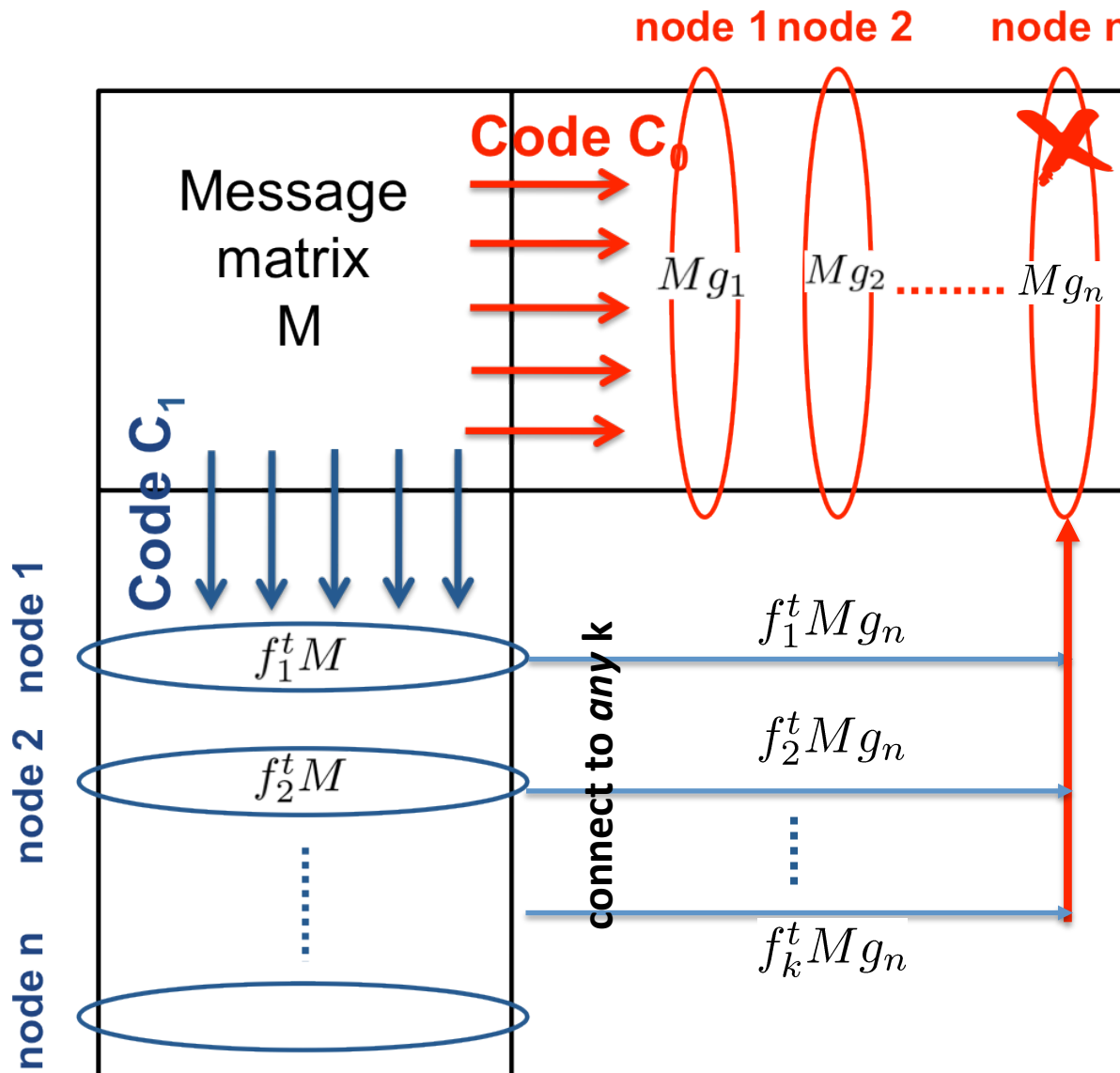
# Twin Code: Data reconstruction



- Connect any k nodes of *one* type
- row-by-row decoding of data

- any k of the $g_i$ are linearly independent

- any k of the $f_i$ are linearly independent

# Twin Code: Repair



- any k of the $g_i$ are linearly independent
- any k of the $f_i$ are linearly independent