# Collective Communication: Theory and Practice

## Robert van de Geijn

# Acknowledgements

This presentation is based on work in the mid-1990s that was sponsored in part by the Intel Research Council and Intel Scalable Systems Division.  At that time, David Payne, Lance Shuler, and Jerrell Watts contributed to the research
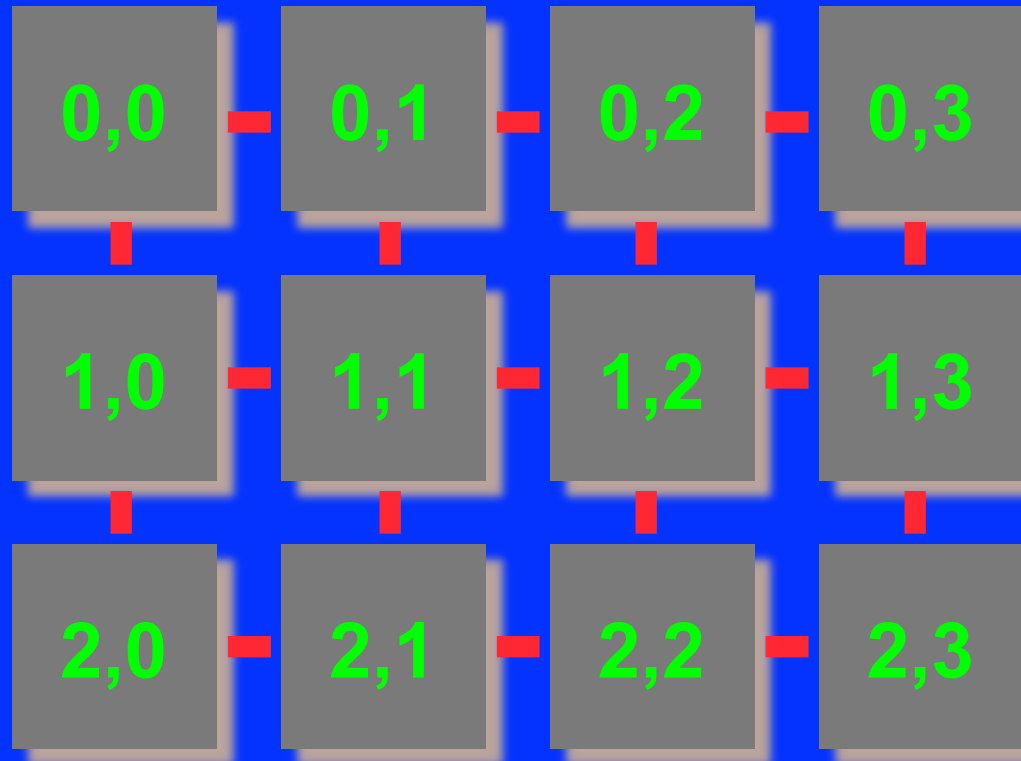
# Outline

Part I:  Theory

- **Model of parallel computation**

- **Collective communications**

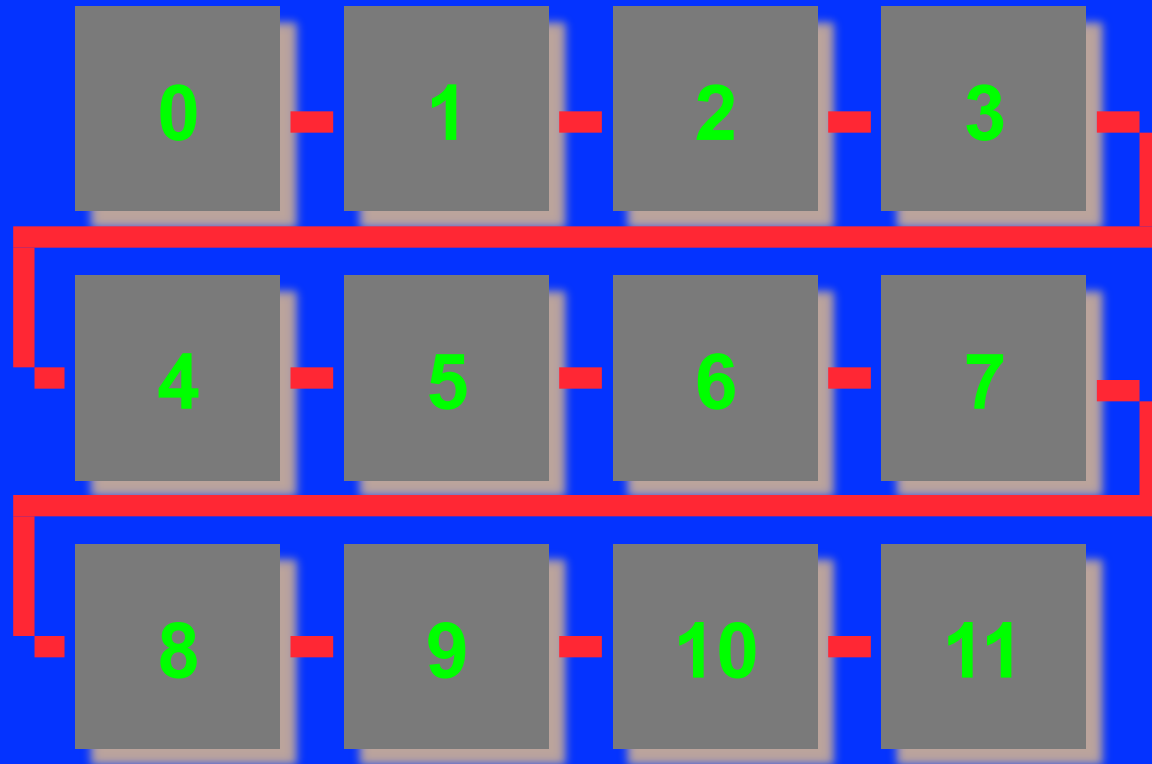- **A building block approach to library implementation**

Part II: Practice

- **Implementation on the Paragon**

- **Performance results**

3

# Outline

Part I:  Theory

- **Model of parallel computation**

- **Collective communications**

- **A building block approach to library implementation**

Part II: Practice

- **Implementation on the Paragon**

- **Performance results**

- **Applications**

4

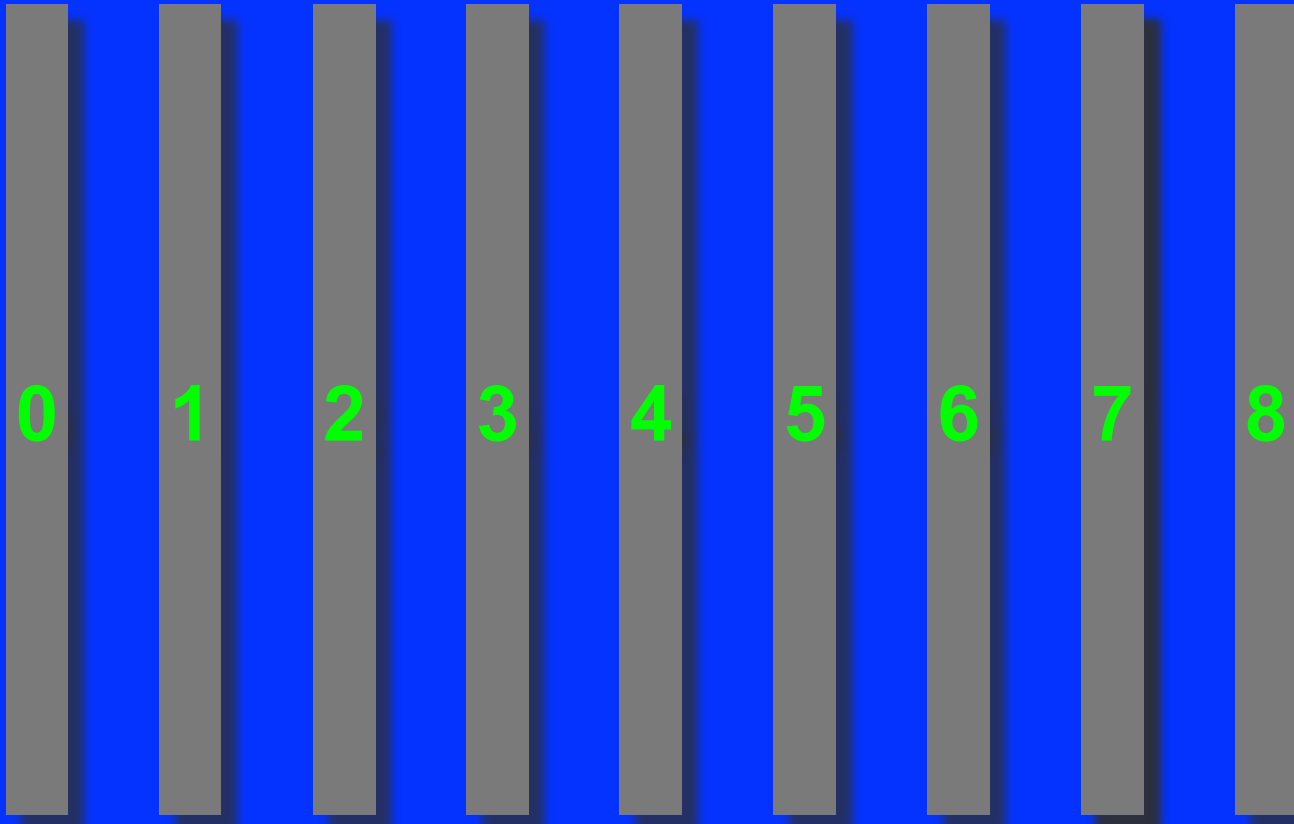# Model of Parallel Computation

- $p$  nodes

- **physical two dimensional mesh**
  - $r$  rows, $c$  columns
  - nodes have physical indices  *(i,j)*

- **often logically viewed as a linear array**
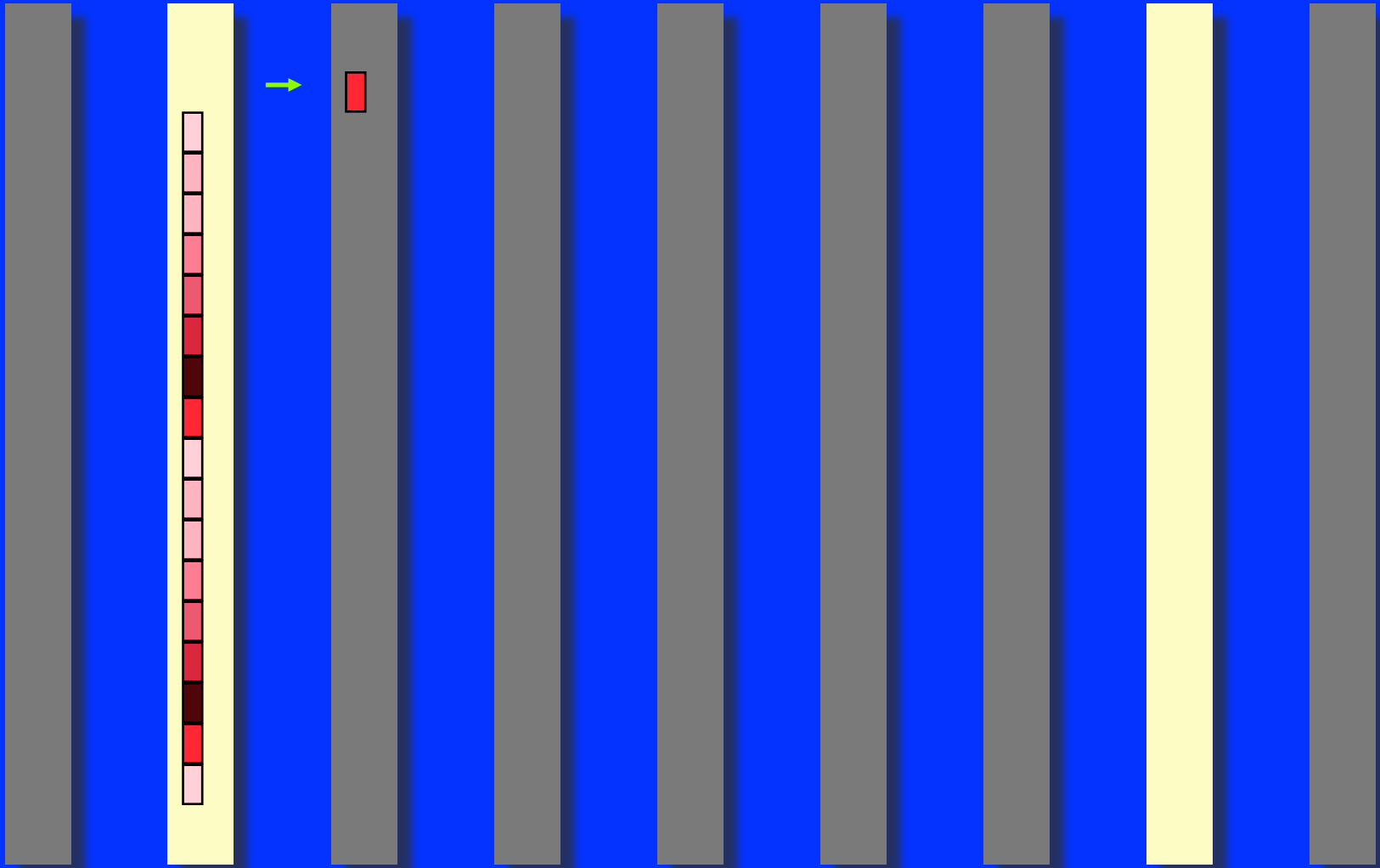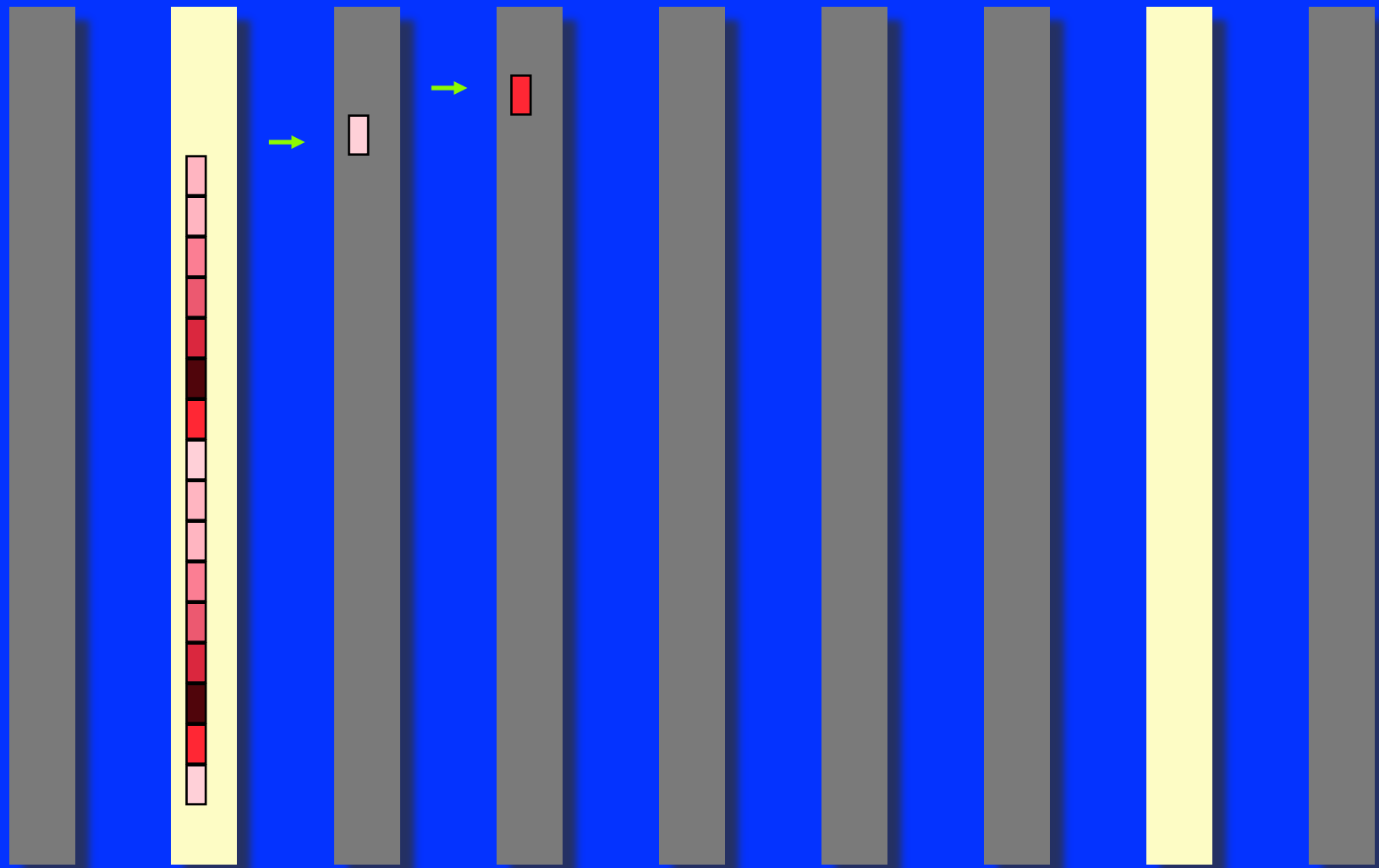  - indexed *0, ... , p-1*
  - nodes are numbered in row-major order

- **physical two dimensional mesh**
  - *r* **rows,** *c* **columns**
  - **nodes have physical indices** *(i,j)*

6

- **often logically viewed as a linear array**
  - indexed *0, ... , p-1*
  - nodes are numbered in row-major order

0 1 2 3 4 5 6 7 8

- **often logically viewed as a linear array**
  - indexed *0, ... , p-1*

8

# The Cost of Communication

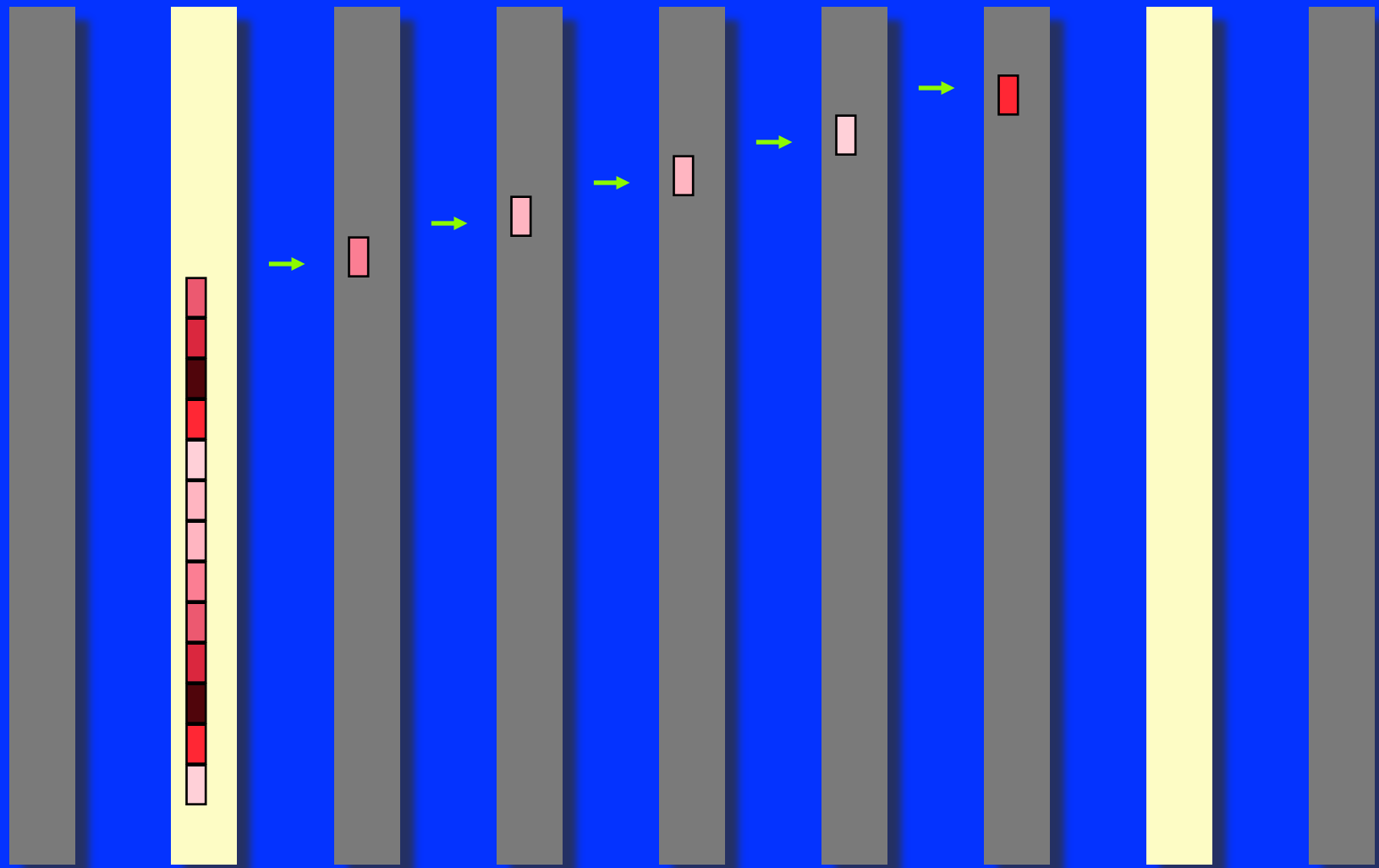- send a message of length $n$ over $d$ links
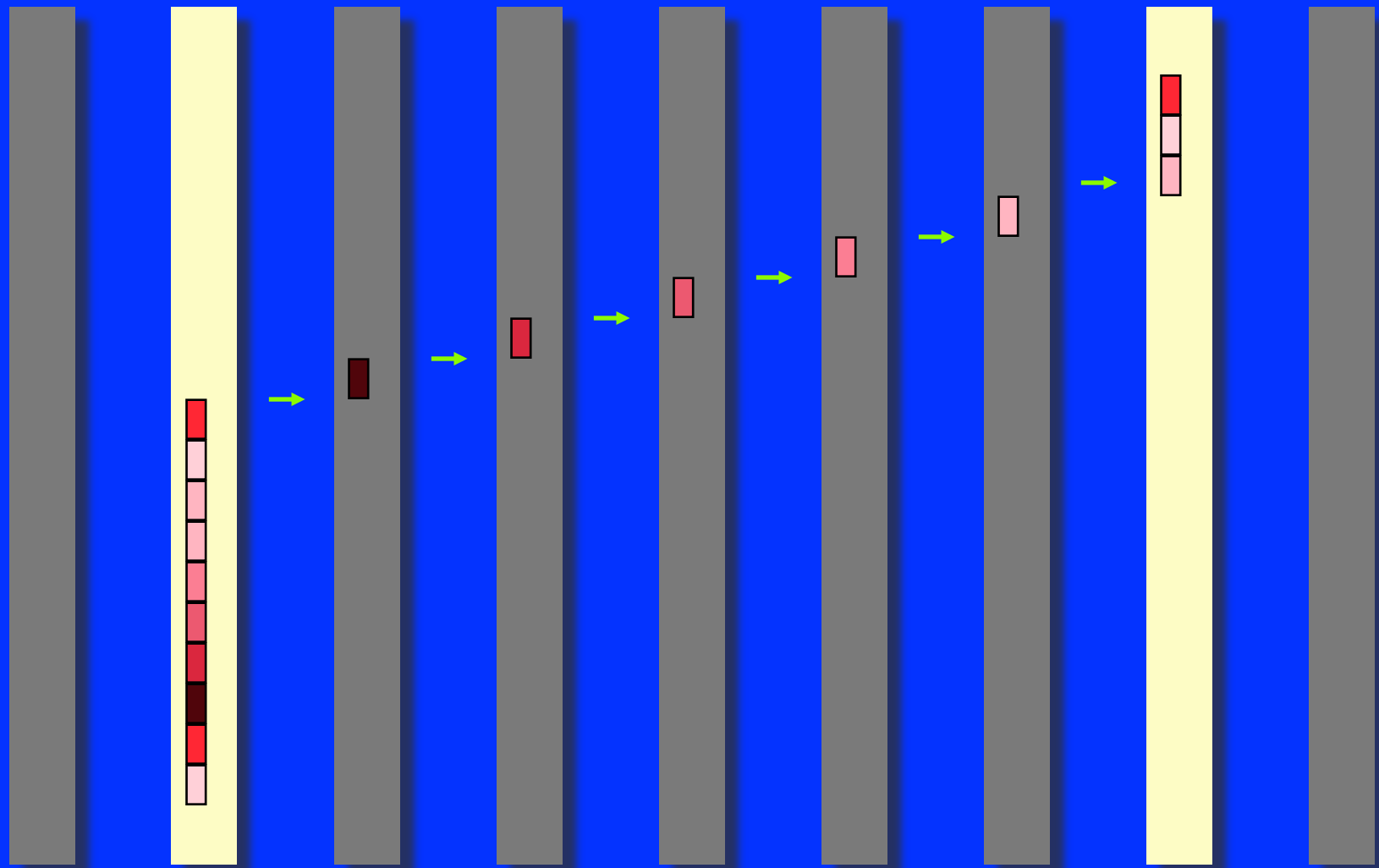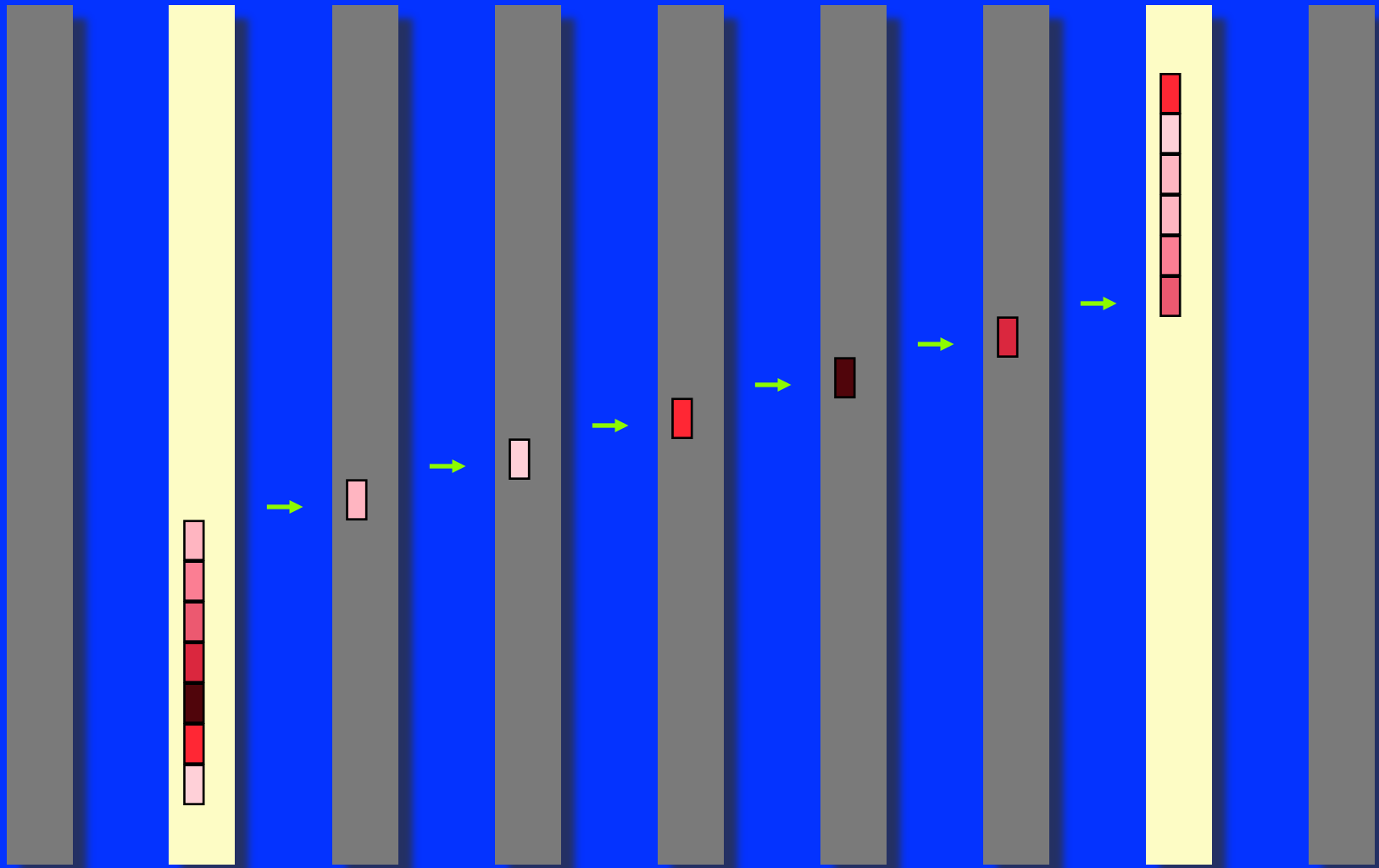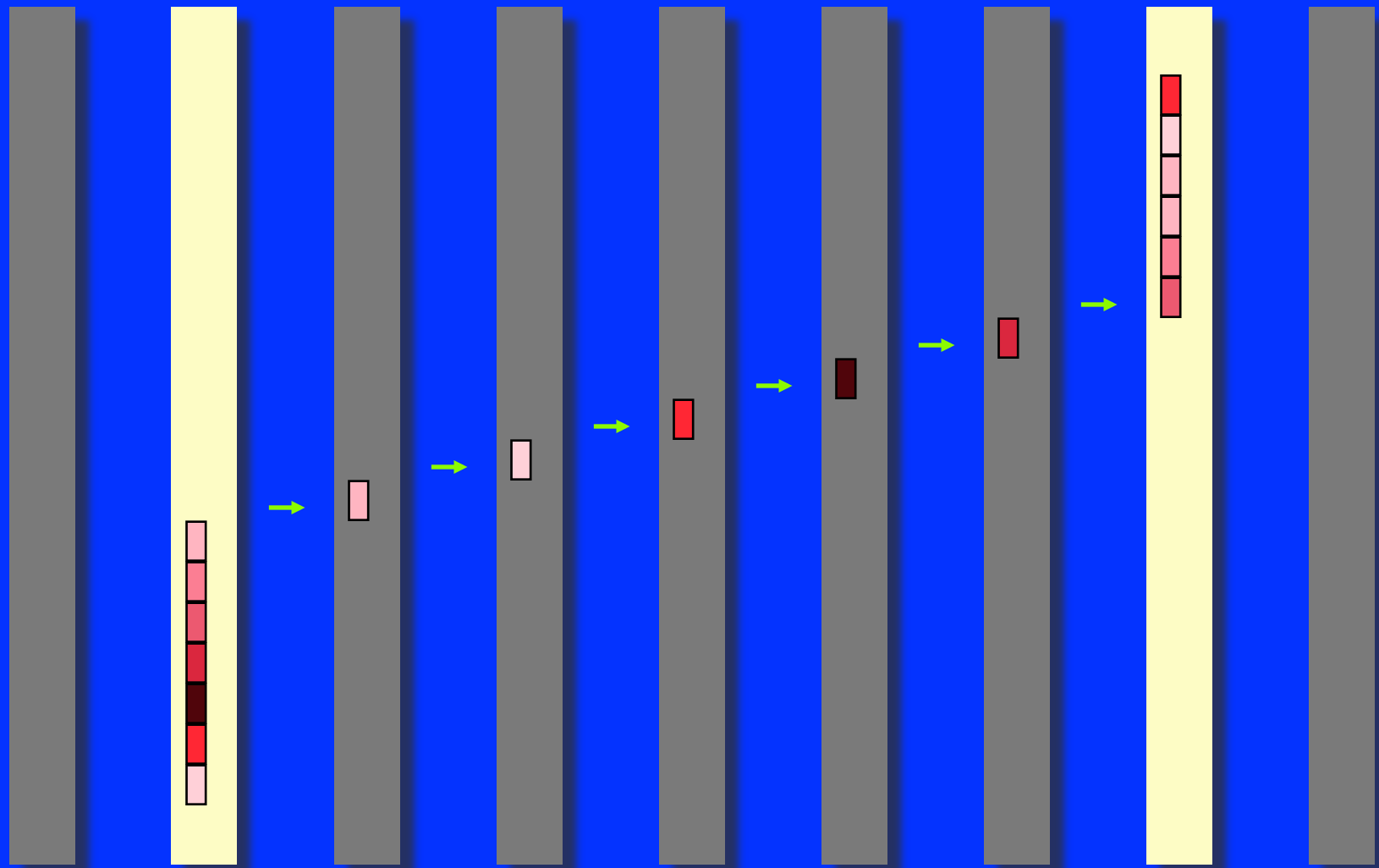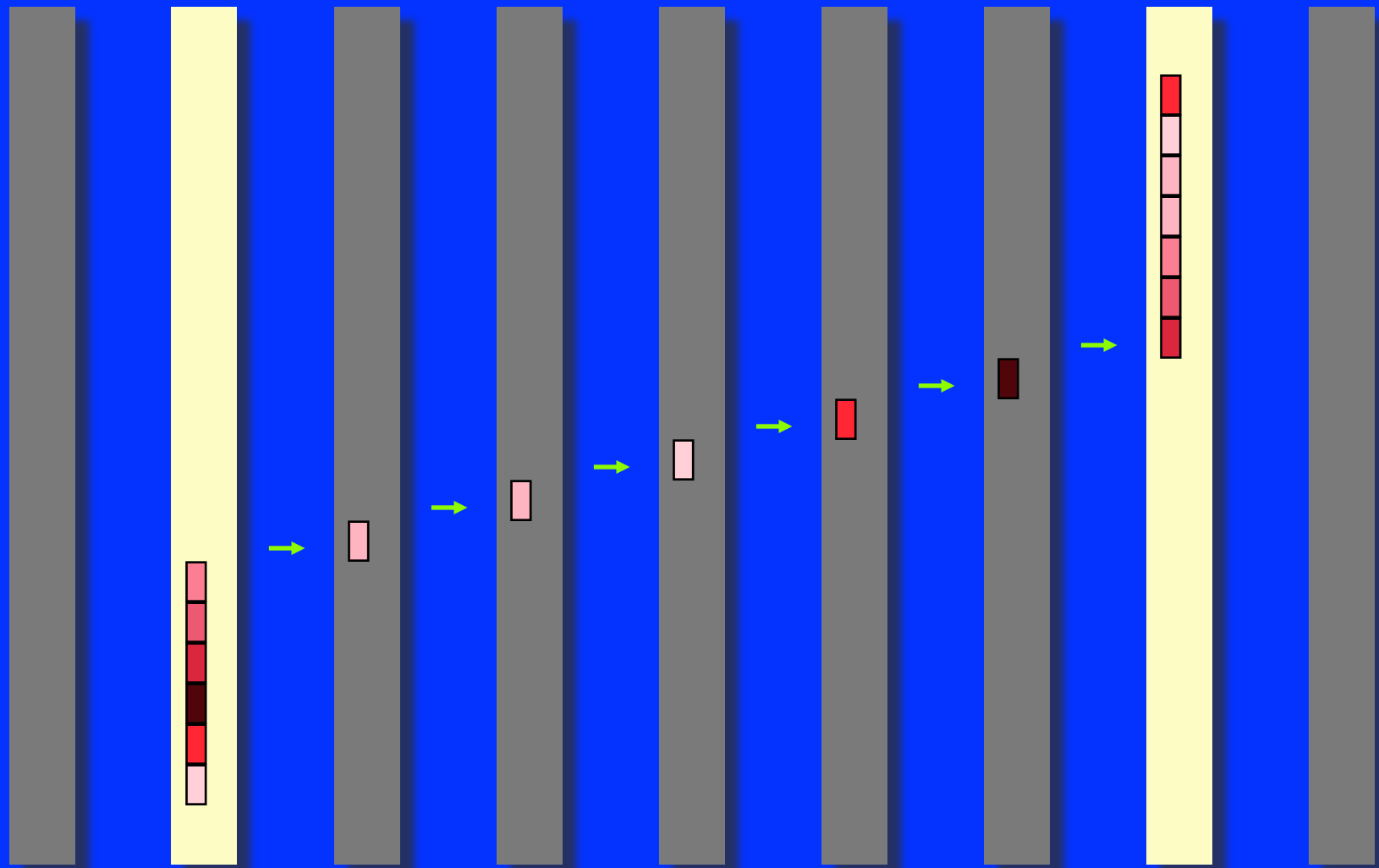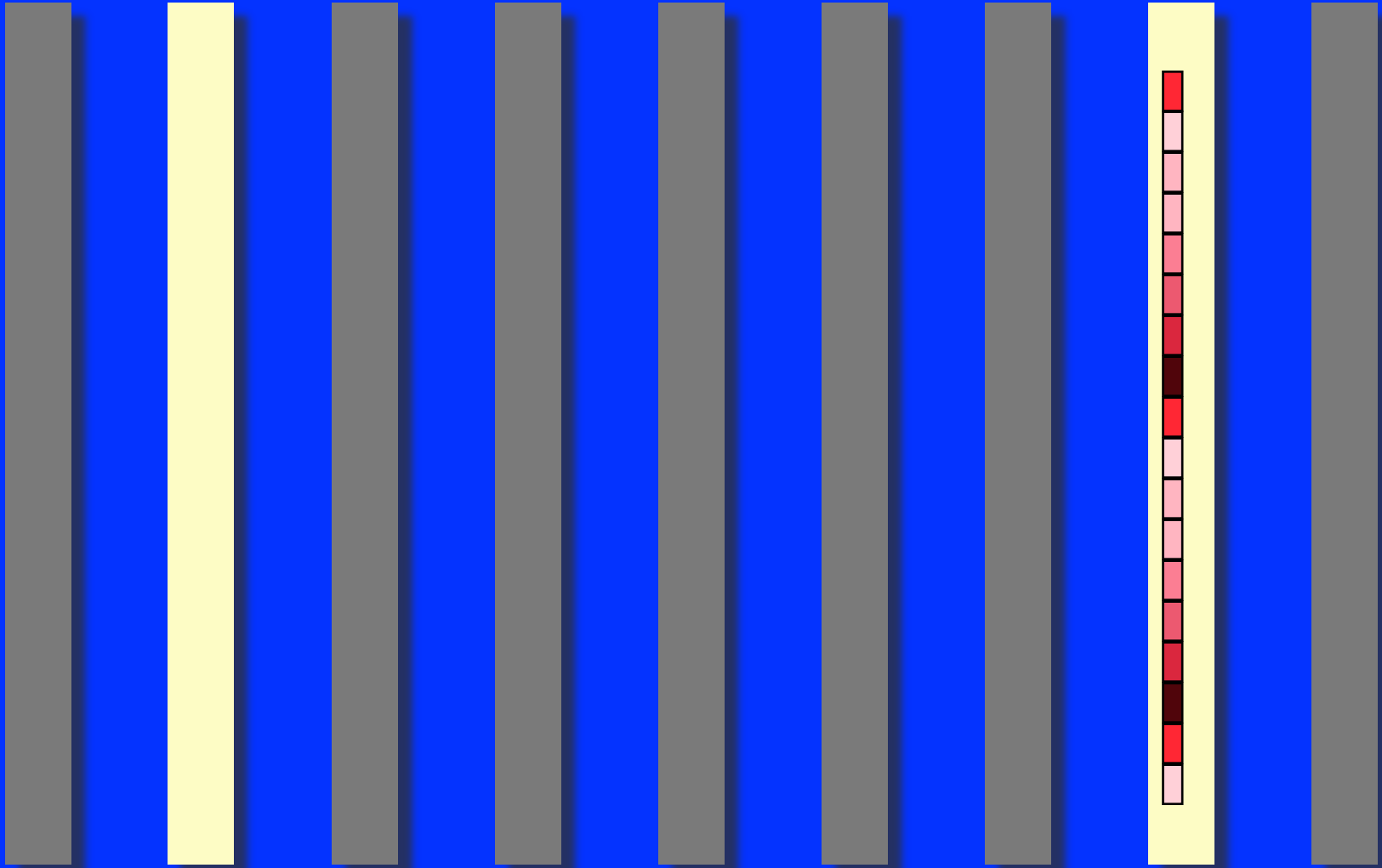- packetize the message
- Example: $d=6$
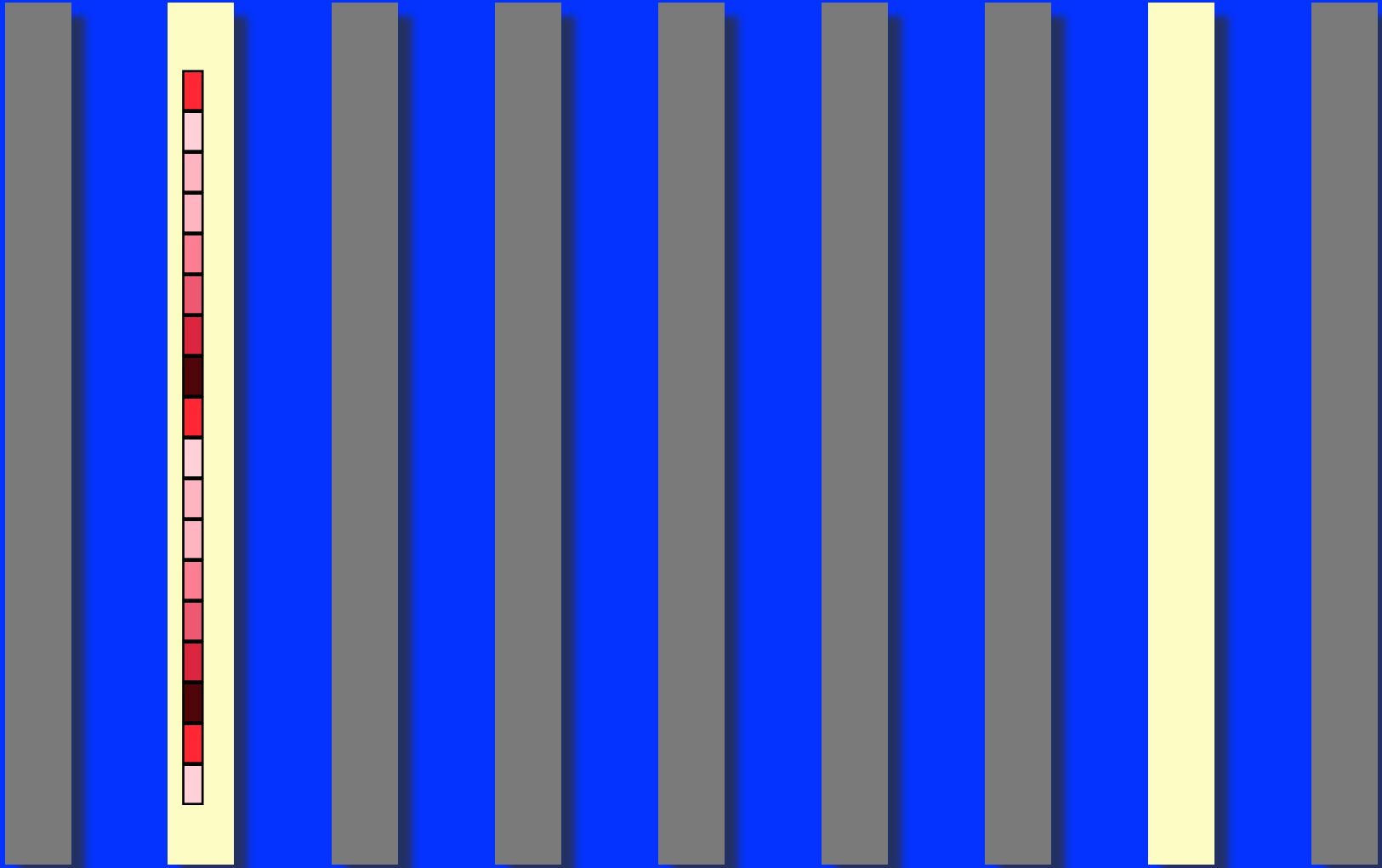
20

# The Cost of Communication

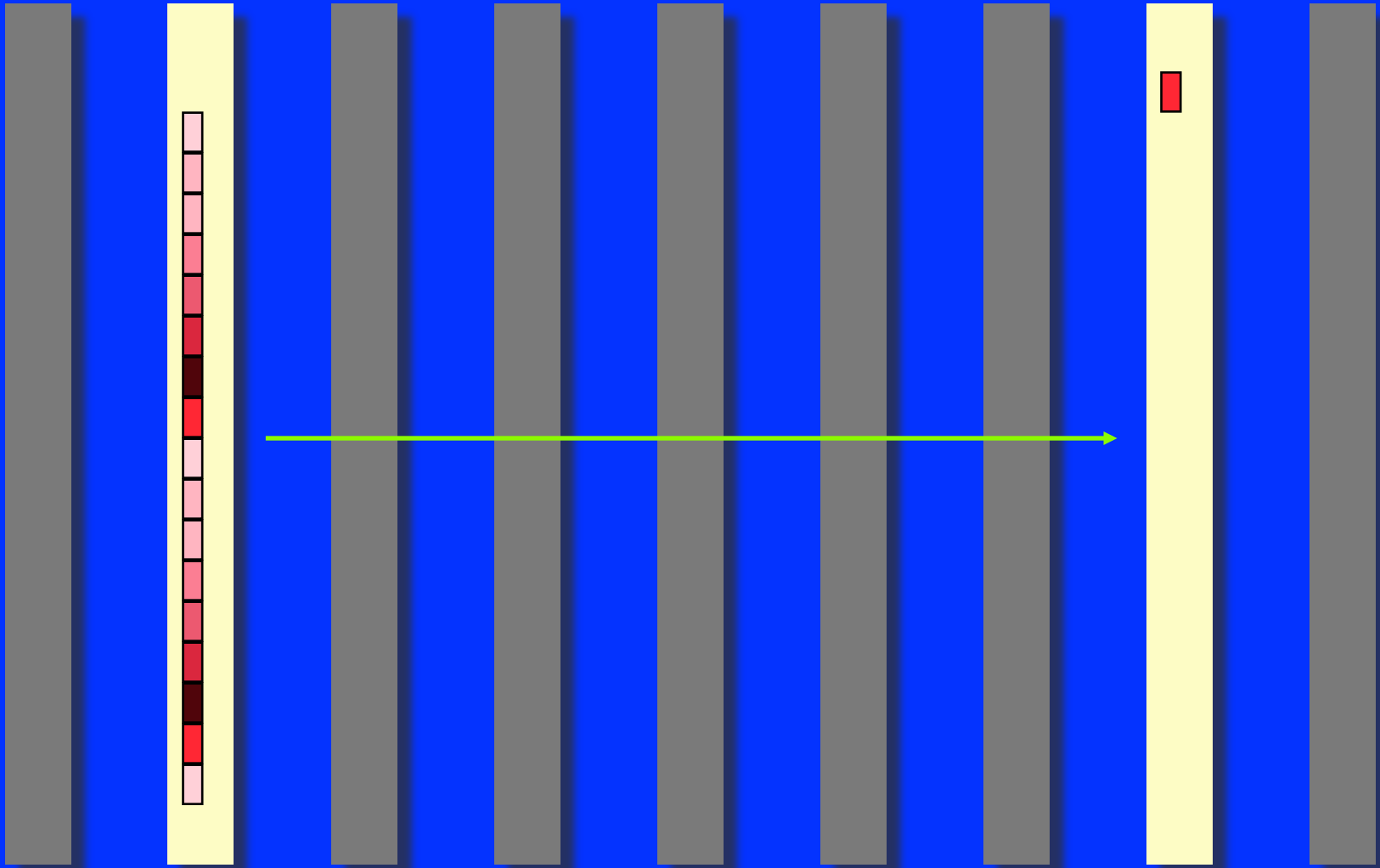- send a message of length $n$ over $d$ links
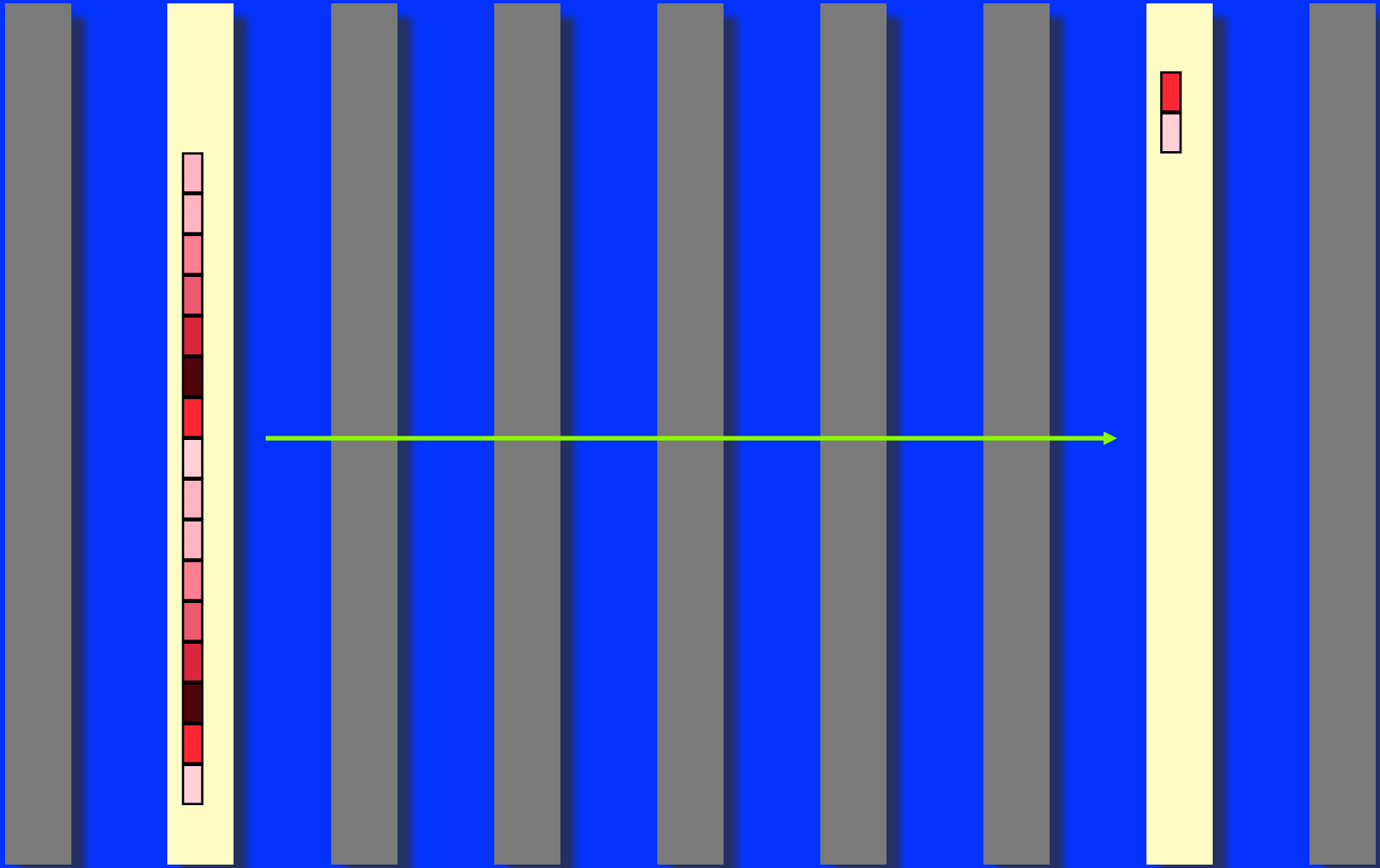- $k$ packets
- Cost:

# The Cost of Communication

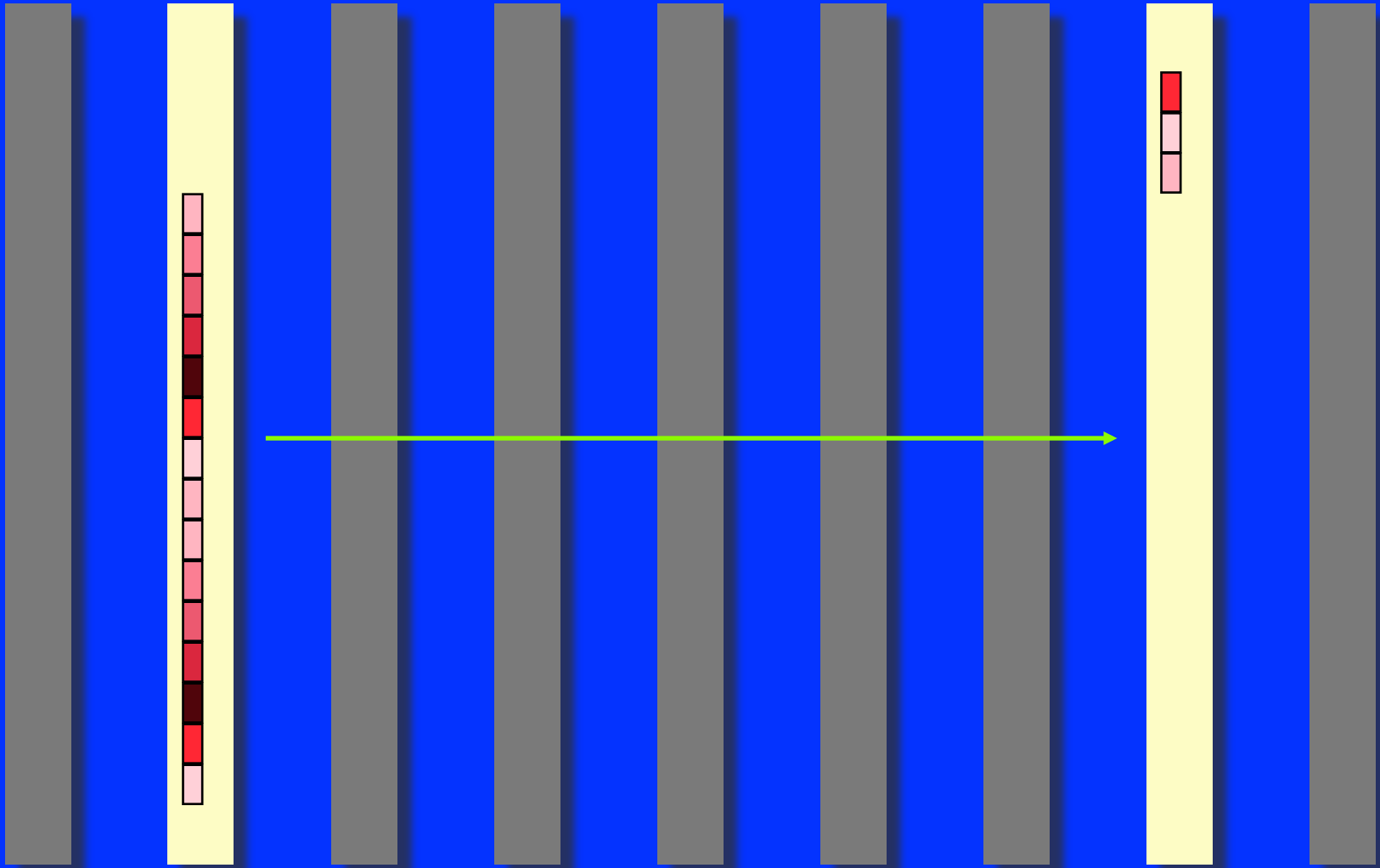- send a message of length $n$ over $d$ links
- $k$ packets
- Cost:

$$\alpha + d\left(\alpha_{net} + \frac{n}{k}\beta\right) + (k-1)\left(\alpha_{net} + \frac{n}{k}\beta\right)$$

$$=$$

$$\alpha + n\beta + (d+k-1)\alpha_{net} + \frac{d-1}{k}\,n\beta$$

$$\approx$$

$$\alpha + n\beta$$

- Example revisited ...

# Model of Parallel Computation

- **a node can send directly to any other node**
- **a node can simultaneously receive and send**
- **cost of communication**
  - sending a message of length $n$ between any two nodes


  - if a message encounters a link that simultaneously accomodates $M$ messages, the cost becomes

# Model of Parallel Computation

- **a node can send directly to any other node**
- **a node can simultaneously receive and send**
- **cost of communication**
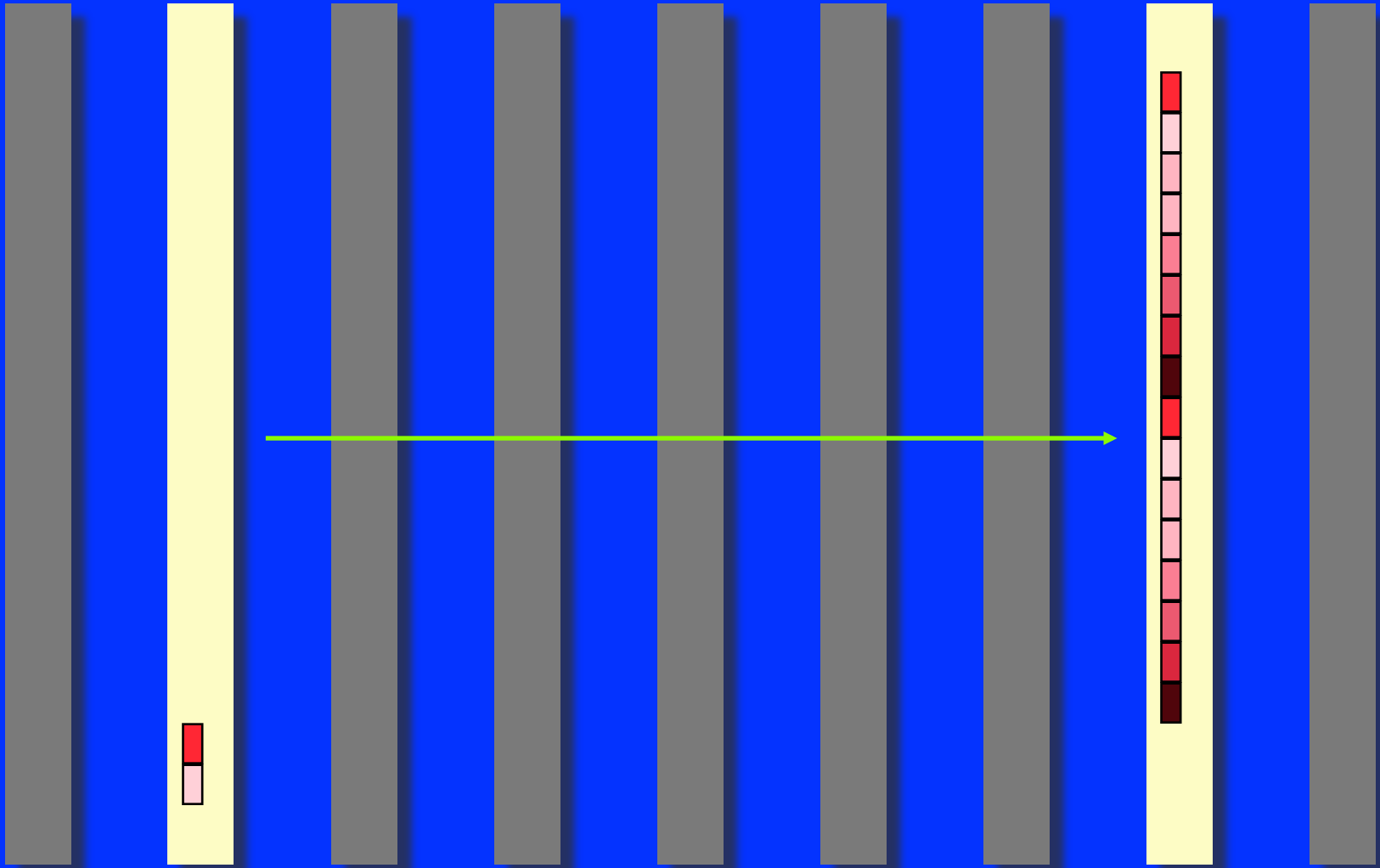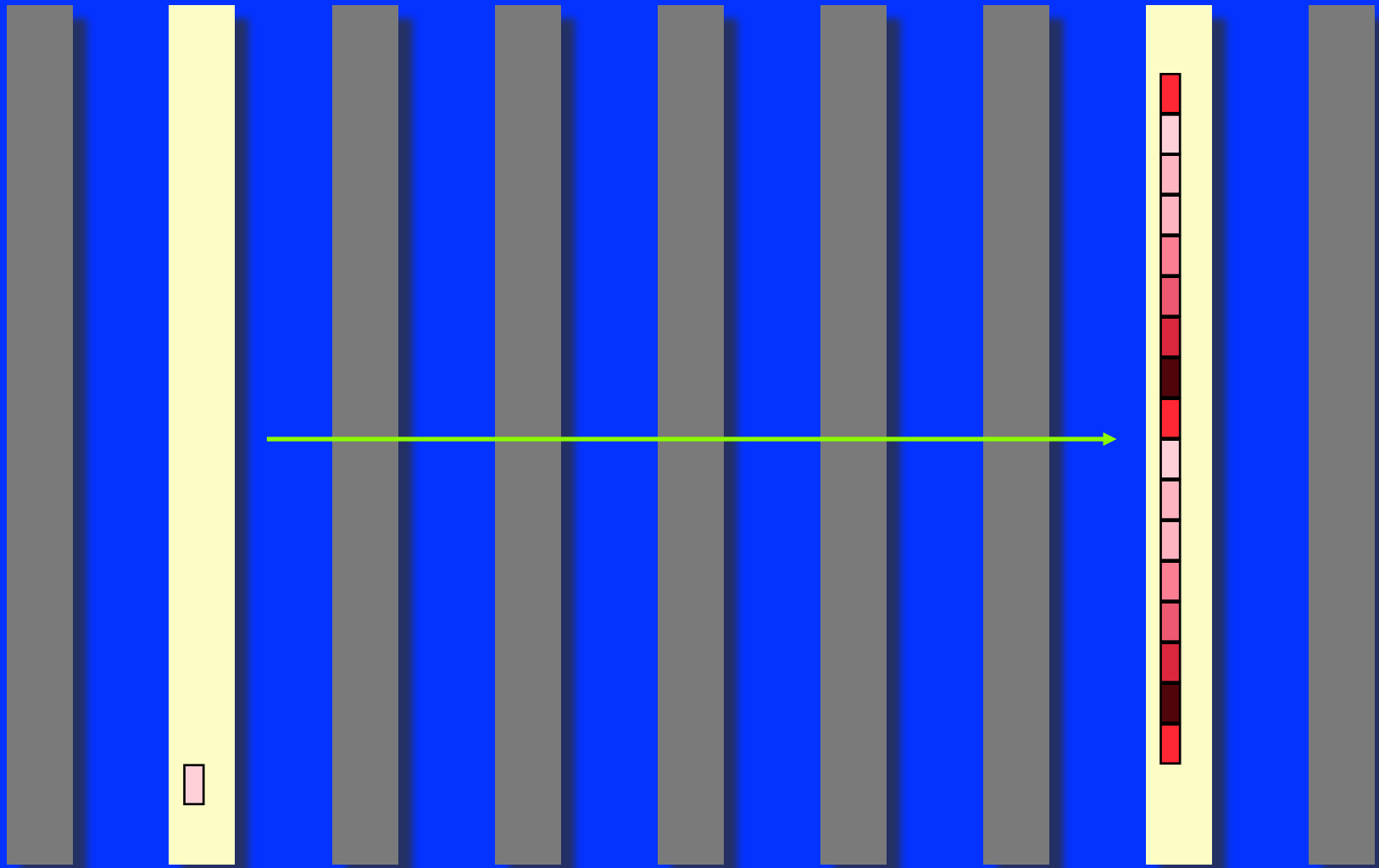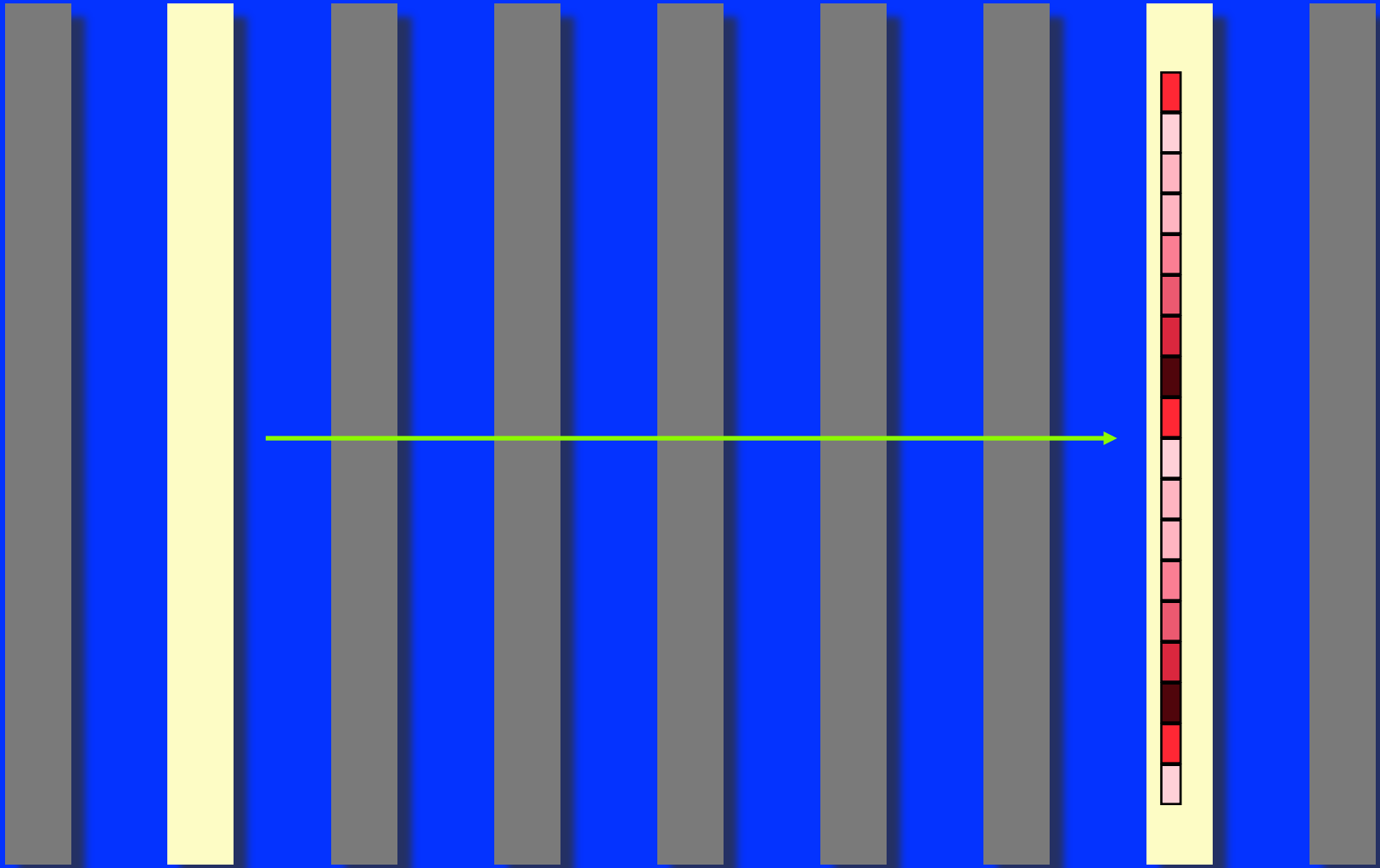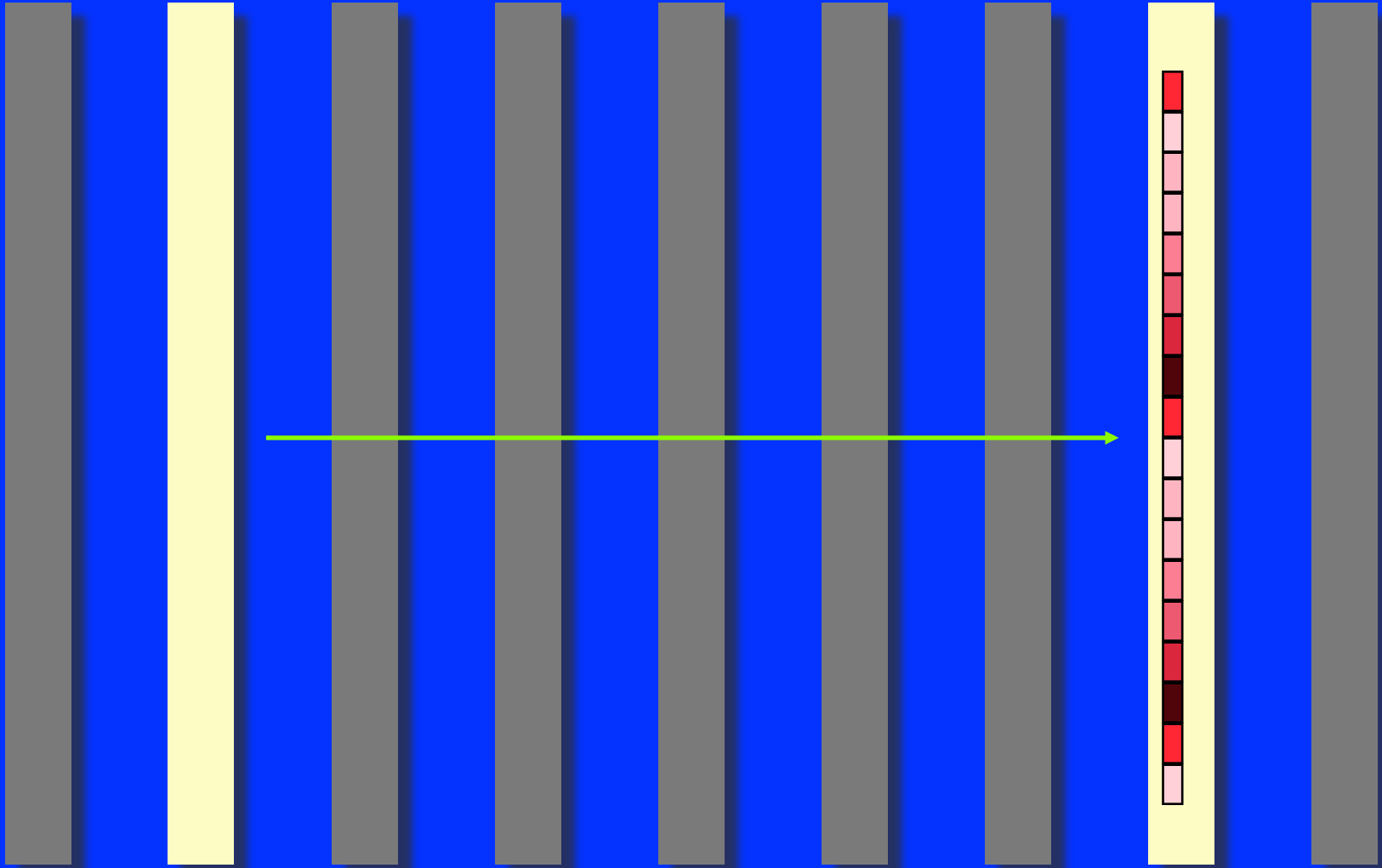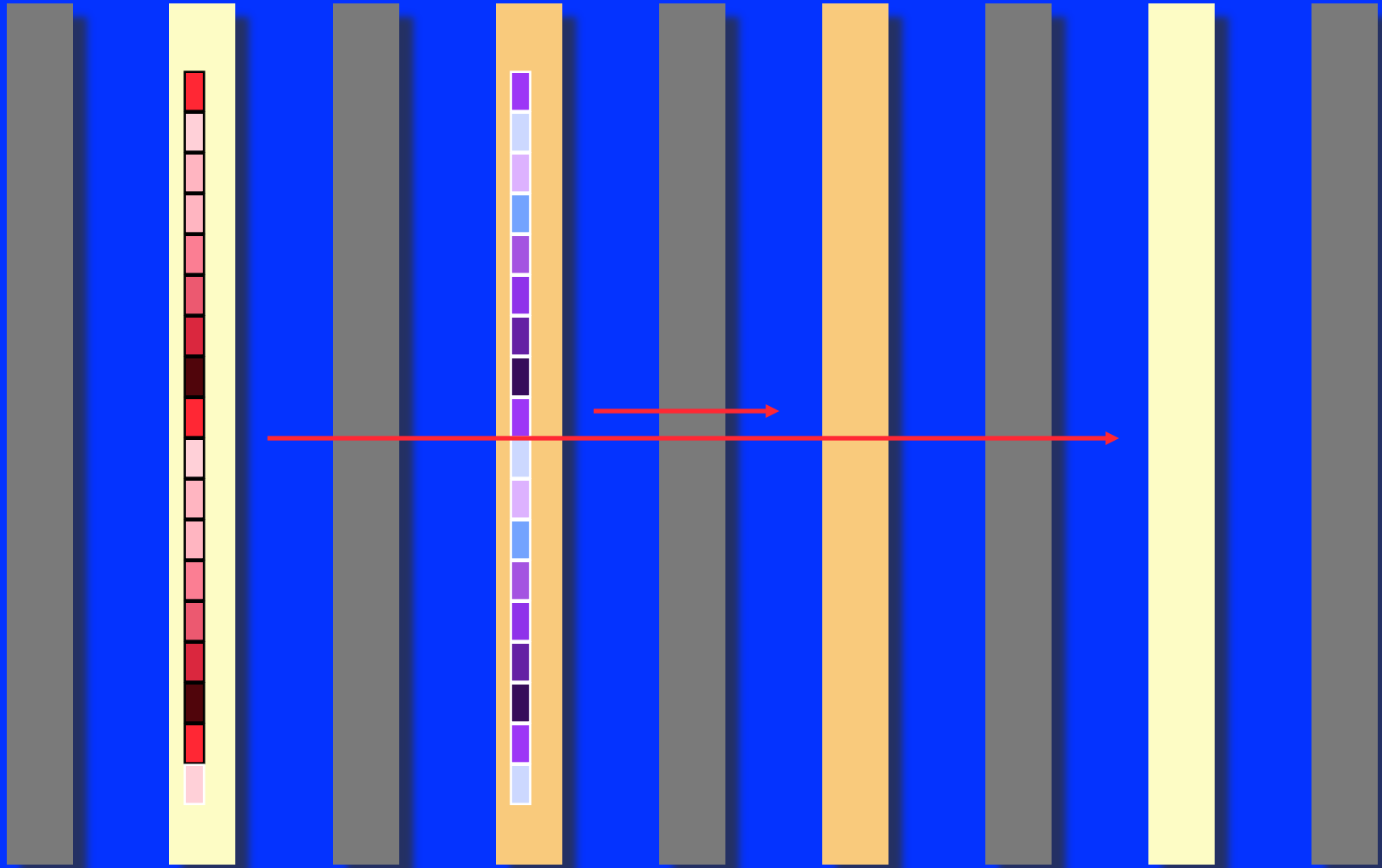  - sending a message of length $n$ between any two nodes
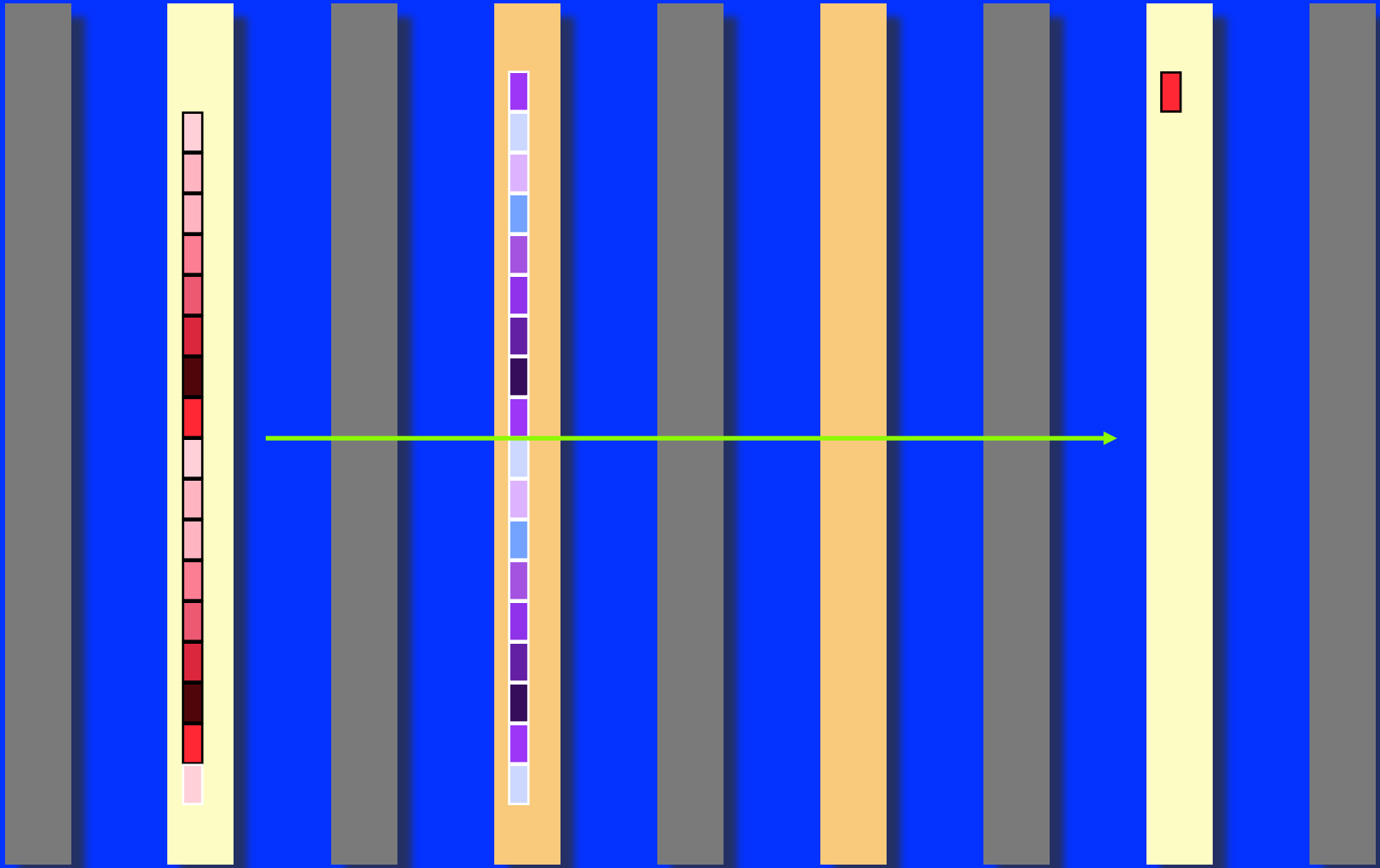
$$\alpha + n\beta$$

  - if a message encounters a link that simultaneously accomodates $M$ messages, the cost becomes
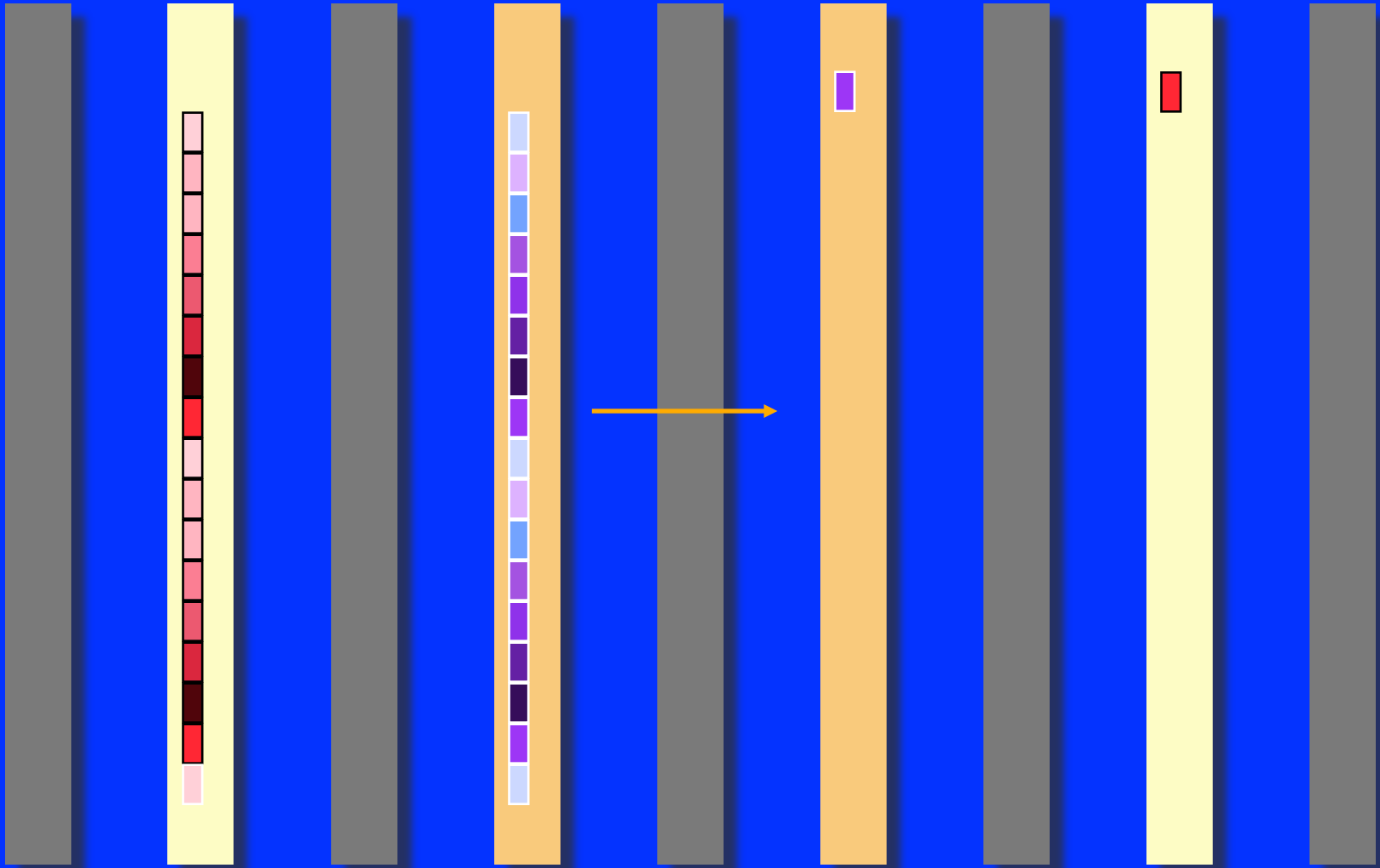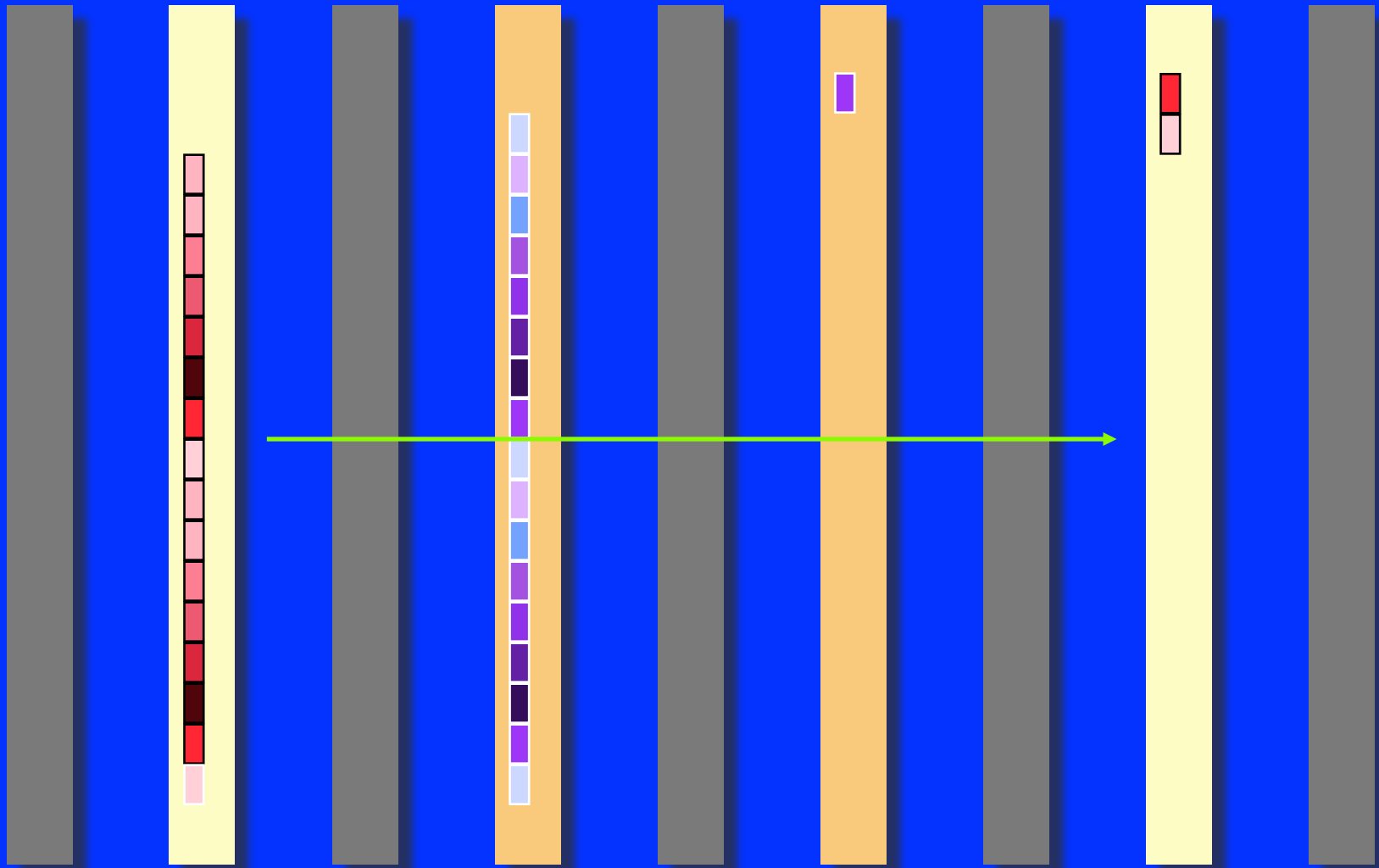
$$\alpha + Mn\beta$$

50

# Interfering messages

- **Example: two messages of length $n$ which share at least one link**

53

63

81

# Outline

**Part I: Theory**

- **Model of computation**

- **Collective communications**

- **A building block approach to library implementation**


**Part II: Practice**

- **Implementation on the Paragon**

- **Performance results**

- **Applications**

91

# Collective Communications

- **Broadcast**
- **Reduce(-to-one)**
- **Scatter**
- **Gather**
- **Allgather**
- **Reduce-scatter**
- **Allreduce**

# Broadcast

**Before**

**After**

93

# Reduce(-to-one)

**Before**

**After**

# Broadcast/Reduce(-to-one)

**Broadcast**

**Reduce(-to-one)**

95

# Scatter

**Before**

**After**

# Gather

**Before**

**After**

97

# Scatter/Gather



Scatter

Gather

# Allgather

**Before**

**After**

99

# Reduce-scatter

**Before**

**After**



100

# Allgather/Reduce-scatter



Allgather

Reduce-scatter

# Allreduce

**Before**                                              **After**

# Lower bounds (startup)

- **Broadcast**     $\lceil log(p) \rceil \alpha$

- **Reduce(-to-one)**     $\lceil log(p) \rceil \alpha$

- **Scatter/Gather**     $\lceil log(p) \rceil \alpha$

- **Allgather**     $\lceil log(p) \rceil \alpha$

- **Reduce-scatter**     $\lceil log(p) \rceil \alpha$

- **Allreduce**     $\lceil log(p) \rceil \alpha$

# Lower bounds (bandwidth)

- **Broadcast**

- **Reduce(-to-one)**

- **Scatter/Gather**

- **Allgather**

- **Reduce-scatter**

- **Allreduce**

$$n\beta$$

$$n\beta + \frac{p-1}{p}n\gamma$$

$$\frac{p-1}{p}n\beta$$

$$\frac{p-1}{p}n\beta$$

$$\frac{p-1}{p}n\beta + \frac{p-1}{p}n\gamma$$

$$2\frac{p-1}{p}n\beta + \frac{p-1}{p}n\gamma$$

# Outline

Part I:  Theory

- **Model of computation**

- **Collective communications**
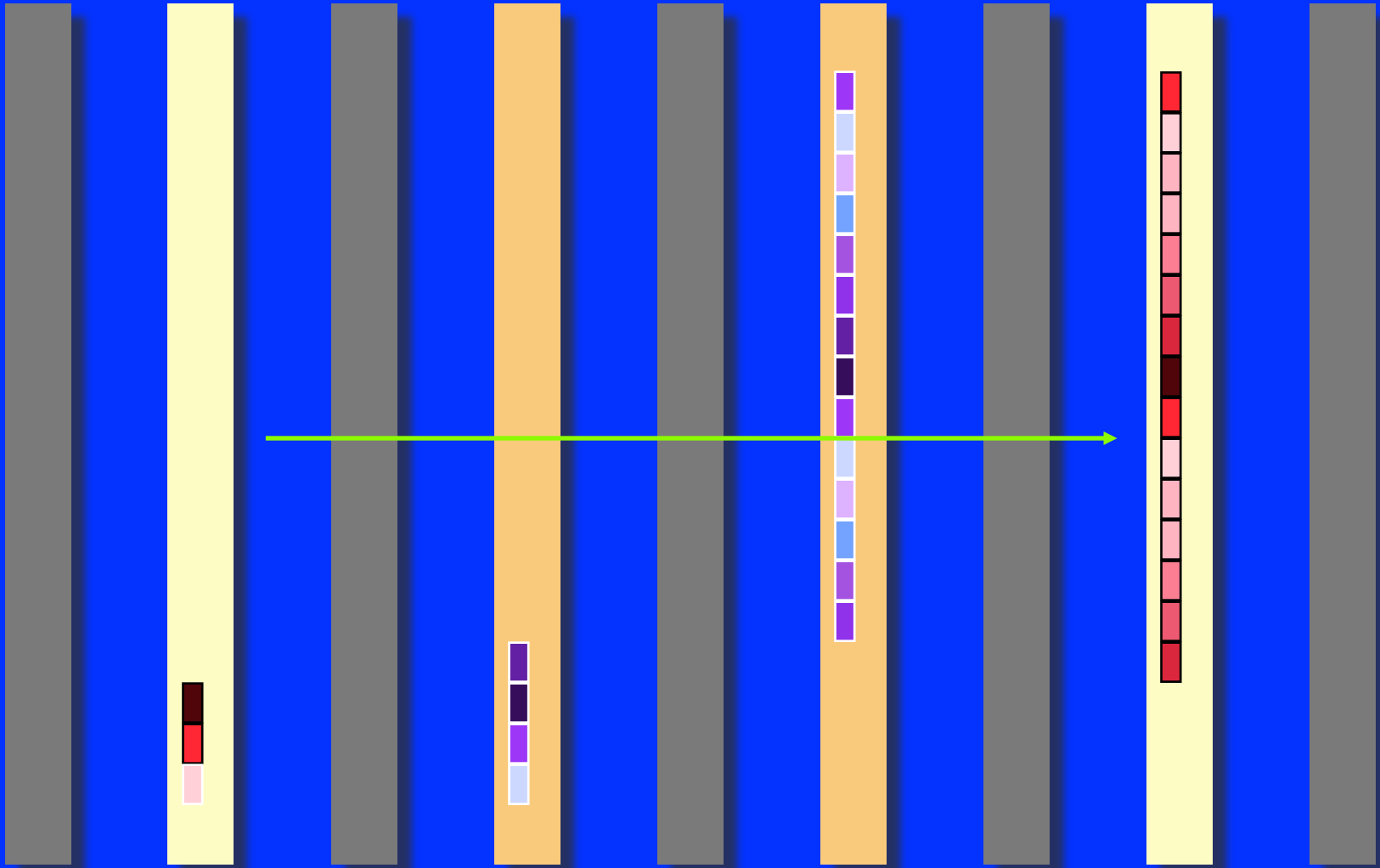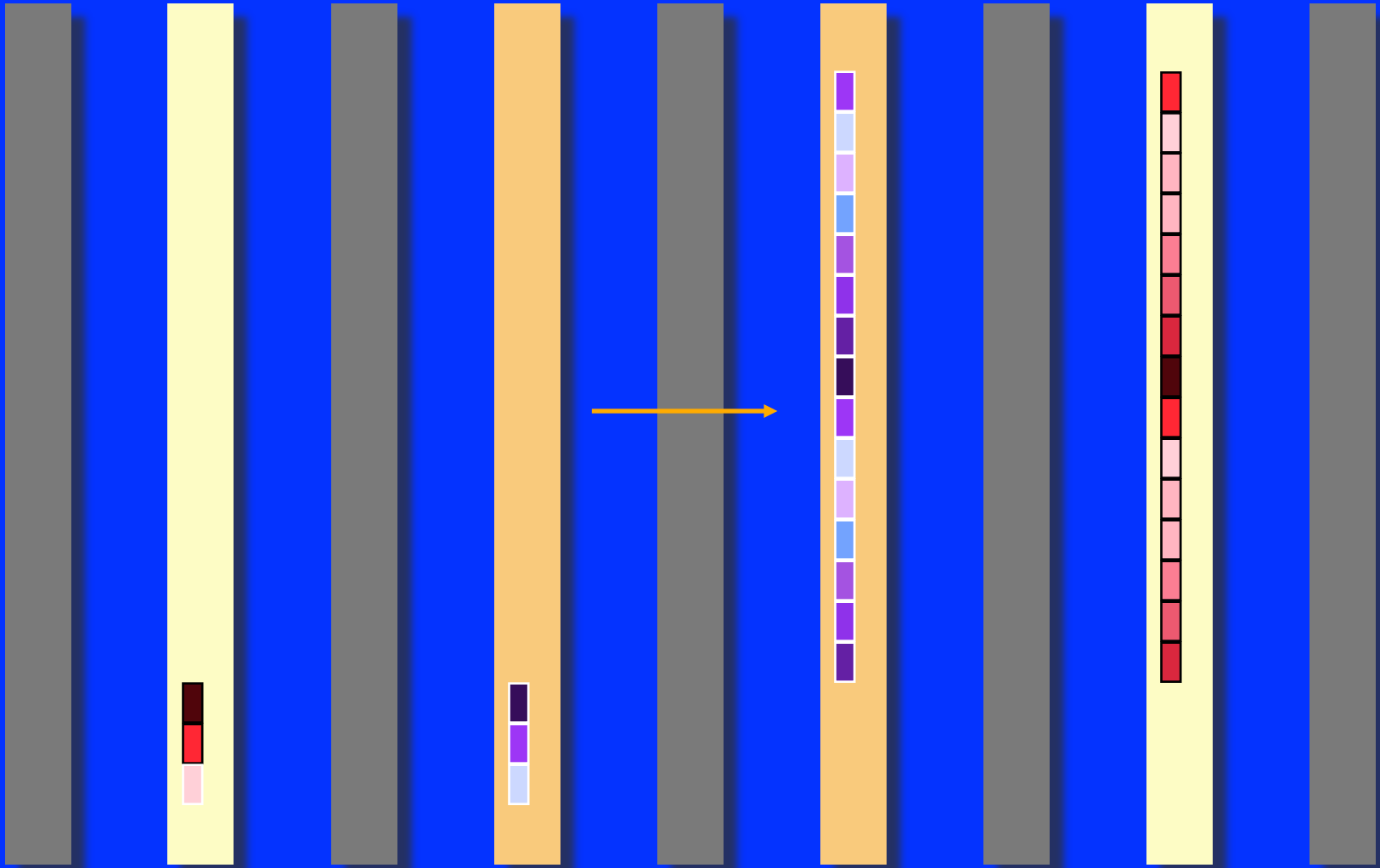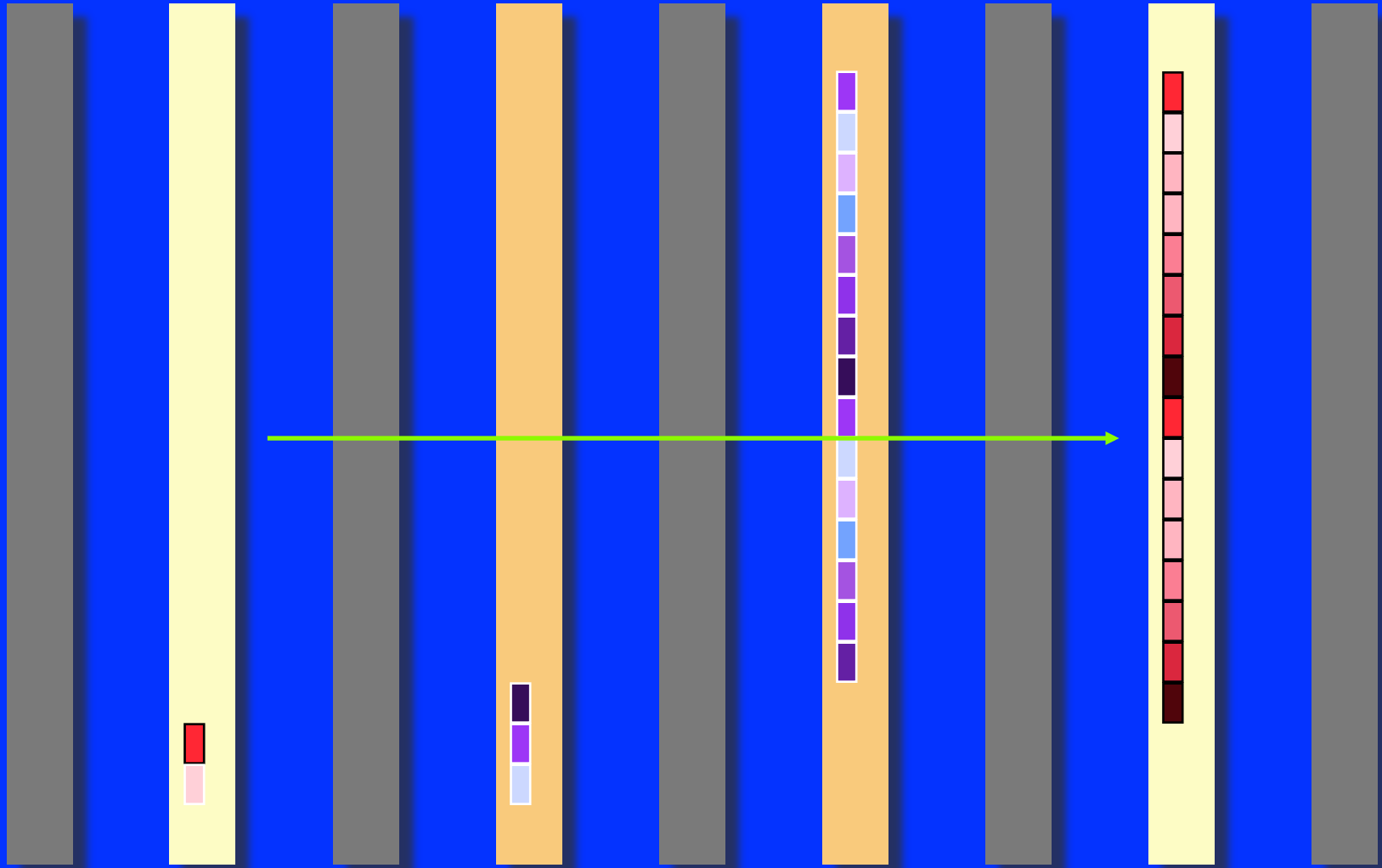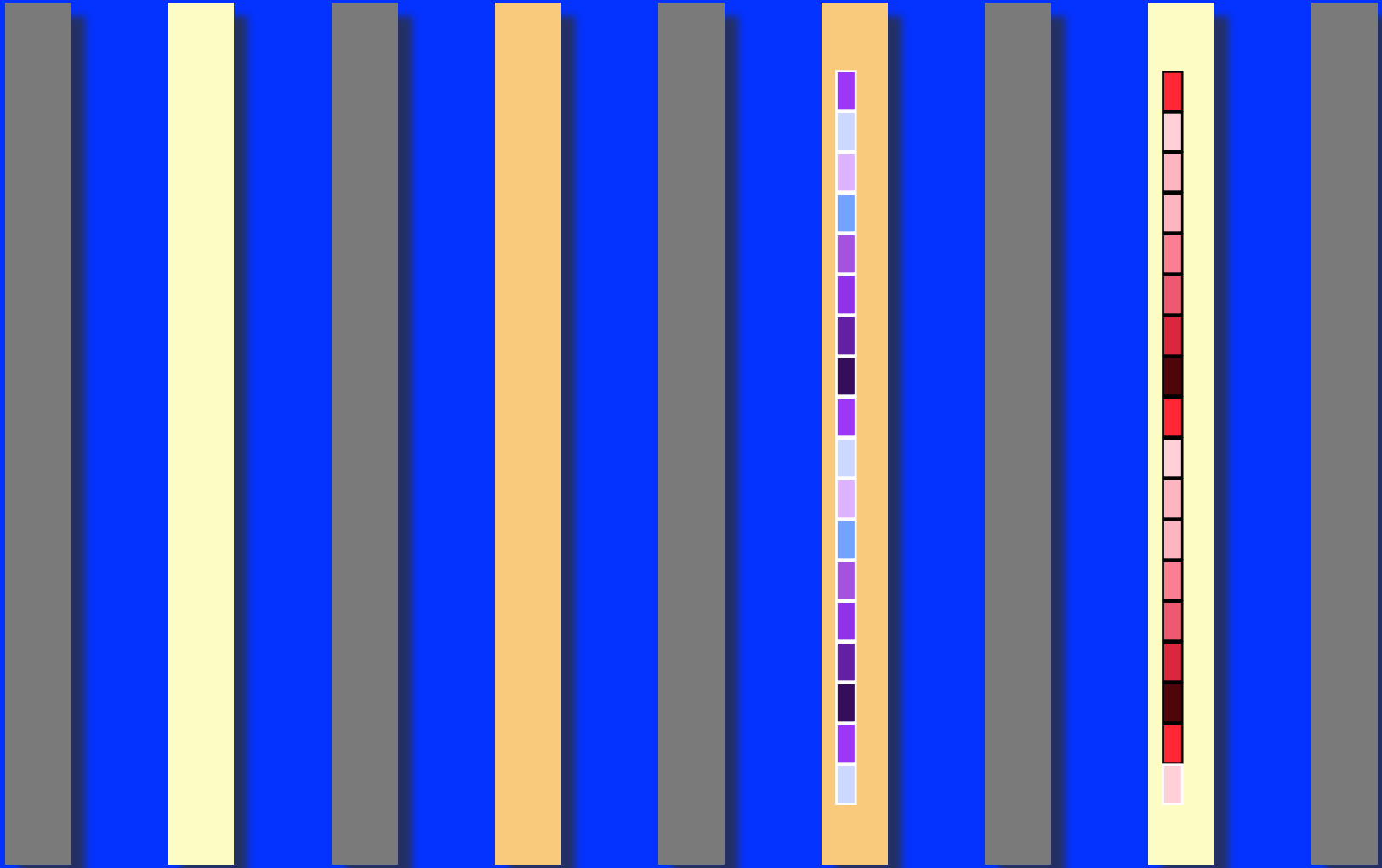
- **A building block approach to library implementation**
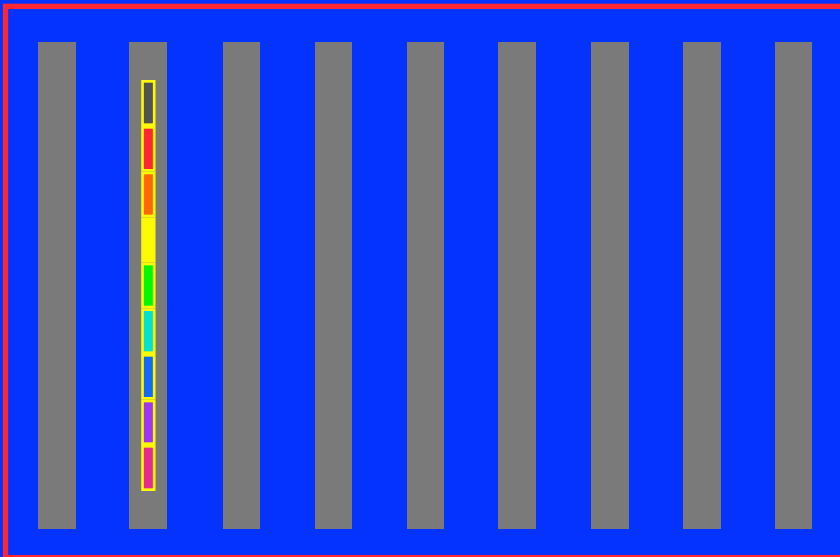
Part II: Practice

- **Implementation on the Paragon**
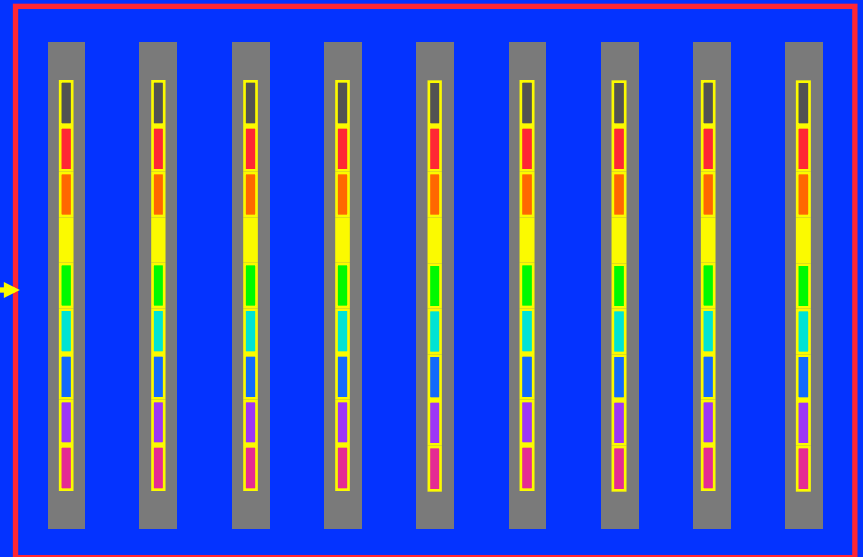
- **Performance results**

- **Applications**

# A building block approach to library implementation

- **Short vector case**

- **Long vector case**

- **Hybrid algorithms**

106

# Short vector case

- **Primary concern:**
  - algorithms must have low latency cost

- **Secondary concerns:**
  - algorithms must work for arbitrary number of nodes
    - » in particular, not just for power-of-two numbers of nodes
  - algorithms should avoid network conflicts
    - » not absolutely necessary, but nice if possible

107

# Minimum spanning tree based algorithms

- **We will show how the following building blocks:**
  - broadcast/combine-to-one
  - scatter/gather

  **can be implemented using minimum spanning trees embedded in the logical linear array while attaining**
  - minimal latency
  - implementation for arbitrary numbers of nodes
  - no network conflicts

108

# General principles

- **message starts on one processor**

# General principles



- **divide logical linear array in half**

110

# General principles

- **send message to the half of the network that does not contain the current node (root) that holds the message**

111

# General principles



- **send message to the half of the network that does not contain the current node (root) that holds the message**

# General principles



- continue recursively in each of the two halves

# Broadcast

**Before**

**After**

116

118

120

121

# Let us view this more closely

- **Red arrows indicate startup of communication (leading to latency, α)**

- **Green arrows indicate packets in transit (leading to a bandwidth related cost proportional to β and the length of the packet)**

125

126

127

128

130

131

132

133

135

139

140

141

146

150

153

155

156

157

158

159

160

161

162

163

164

# Cost of minimum spanning tree broadcast

$$\lceil log(p) \rceil \ (\alpha + n\beta)$$

number of steps                    cost per steps

# Cost of minimum spanning tree broadcast

$$\lceil log(p) \rceil \ (\alpha + n\beta)$$

number of steps

cost per steps

**Notice: attains lower bound for latency component**

```
MSTBCAST( x, root, left, right )

if left = right return
mid = ⌊(left + right)/2⌋
if root ≤ mid then dest = right else dest = left

if me == root SEND( x, dest )
if me == dest RECV( x, root )

if me ≤ mid and root ≤ mid
    MSTBCAST( x, root, left, mid )
else if me ≤ mid and root > mid
    MSTBCAST( x, dest, left, mid )
else if me > mid and root ≤ mid
    MSTBCAST( x, dest, mid+1, right )
else if me > mid and root > mid
    MSTBCAST( x, root, mid+1, right )
```

left

root

me

dest

mid

right

169

dest

root          dest          me          root          dest          root          dest          root

left          mid          right          left          right          left          right          left          right
                                          mid                        mid                        mid

171

# Reduce(-to-one)

**Before**

**After**

172

174

179

180

181

# Cost of minimum spanning tree reduce(-to-one)

$$\lceil log(p) \rceil \; (\alpha + n\beta + n\gamma)$$

number of steps

cost per steps

# Cost of minimum spanning tree
## reduce(-to-one)

$$\lceil log(p) \rceil \left( \alpha + n\beta + n\gamma \right)$$

**number of steps**

**cost per steps**

**Notice: attains lower bound for latency component**

```
MSTReduce( x, root, left, right )

if left = right return
mid = ⌊(left + right)/2⌋
if root ≤ mid then srce = right else srce = left

if me ≤ mid and root ≤ mid
    MSTReduce( x, root, left, mid )
else if me ≤ mid and root > mid
    MSTReduce( x, srce, left, mid )
else if me > mid and root ≤ mid
    MSTReduce( x, srce, mid+1, right )
else if me > mid and root > mid
    MSTReduce( x, root, mid+1, right )

if me == srce Send( x, root )
if me == root Recv( tmp, srce ) and x = x + tmp
```

# Scatter

**Before**

**After**



185

187

188

189

191

192

193

195

197

198

# Cost of minimum spanning tree scatter

- **Assumption: power of two number of nodes**

$$\sum_{k=1}^{log(p)} \left( \alpha + \frac{n}{2^k} \beta \right)$$
$$=$$
$$log(p) \ \ \alpha + \frac{p-1}{p} n\beta$$

# Cost of minimum spanning tree scatter

- **Assumption: power of two number of nodes**

$$\sum_{k=1}^{log(p)} \left( \alpha + \frac{n}{2^k} \beta \right) = log(p) \ \alpha + \frac{p-1}{p} n\beta$$

**Notice: attains lower bound for latency and bandwidth components**

# Gather

**Before**

**After**



202

204

205

206

208

209

210

211

212

213

214

215

216

# Cost of minimum spanning tree gather

- Assumption: power of two number of nodes

$$\sum_{k=1}^{log(p)} \left( \alpha + \frac{n}{2^k} \beta \right)$$
$$=$$
$$log(p) \ \alpha + \frac{p-1}{p} n\beta$$

# Cost of minimum spanning tree gather

- **Assumption: power of two number of nodes**

$$\sum_{k=1}^{log(p)} \left( \alpha + \frac{n}{2^k} \beta \right)$$
$$=$$
$$log(p) \ \ \alpha + \frac{p-1}{p} n\beta$$

**Notice: attains lower bound for latency and bandwidth components**

# Using the building blocks

# Allgather (short vector)

# Allgather (short vector)



Gather

# Allgather (short vector)

Broadcast

222

# Cost of gather/broadcast allgather

- **Assumption: power of two number of nodes**

**gather**

$$log(p)\alpha + \frac{p-1}{p}n\beta$$

**broadcast**

$$\frac{log(p)(\alpha + n\beta)}{2log(p)\alpha + \left(\frac{p-1}{p} + log(p)\right)n\beta}$$

# Cost of gather/broadcast allgather

- Assumption: power of two number of nodes

**gather**

**broadcast**

$$log(p)\alpha + \frac{p-1}{p}n\beta$$

$$log(p)(\alpha + n\beta)$$

$$2log(p)\alpha + \left(\frac{p-1}{p} + log(p)\right)n\beta$$

Notice: does not attain lower bound for latency or bandwidth components

224

# Reduce-scatter
# (short vector)

# Reduce-scatter
# (short vector)

**Reduce(-to-one)**

226

# Reduce-scatter
# (short vector)

**Scatter**

227

# Cost of Reduce(-to-one)/scatter Reduce-scatter

- **Assumption: power of two number of nodes**

**Reduce(-to-one)** $log(p)(\alpha + n\beta + n\gamma)$

**scatter** $log(p)\alpha + \dfrac{p-1}{p}n\beta$

$$\dfrac{log(p)\alpha + \dfrac{p-1}{p}n\beta}{2log(p)\alpha + \left(\dfrac{p-1}{p} + log(p)\right)n\beta + log(p)n\gamma}$$

228

# Cost of Reduce(-to-one)/scatter reduce-scatter

- **Assumption: power of two number of nodes**

**Reduce(-to-one)** $log(p)(\alpha + n\beta + n\gamma)$

**scatter** $log(p)\alpha + \dfrac{p-1}{p}n\beta$

$$2log(p)\alpha + \left(\dfrac{p-1}{p} + log(p)\right)n\beta + log(p)n\gamma$$

**Notice: does not attain lower bound for latency or bandwidth components**

229

# Allreduce
# (short vector)

# Allreduce
# (short vector)

**Reduce(-to-one)**

231

# Allreduce
# (short vector)



Broadcast

232

# Cost of reduce(-to-one)/broadcast Allreduce

- **Assumption: power of two number of nodes**

**Reduce(-to-one)** $log(p)(\alpha + n\beta + n\gamma)$

**broadcast** $log(p)(\alpha + n\beta)$

$$2log(p)\alpha + 2log(p)n\beta + log(p)n\gamma$$

233

# Cost of reduce(-to-one)/broadcast Allreduce

- **Assumption: power of two number of nodes**

**Reduce(-to-one)** $log(p)(\alpha + n\beta + n\gamma)$

**broadcast** $log(p)(\alpha + n\beta)$

$$2log(p)\alpha + 2log(p)n\beta + log(p)n\gamma$$

Notice: does not attain lower bound for latency or bandwidth components

# Recap

**Reduce(-to-one)**

$log(p)(\alpha + n\beta + n\gamma)$

**Scatter**

$log(p)\alpha + \dfrac{p-1}{p}n\beta$

**Gather**

$log(p)\alpha + \dfrac{p-1}{p}n\beta$

**Broadcast**

$log(p)(\alpha + n\beta)$

**Reduce-scatter**

**Allreduce**

**Allgather**

235

# Recap

**Reduce(-to-one)**

$log(p)(\alpha + n\beta + n\gamma)$

**Scatter**

$log(p)\alpha + \dfrac{p-1}{p}n\beta$

**Gather**

$log(p)\alpha + \dfrac{p-1}{p}n\beta$

**Broadcast**

$log(p)(\alpha + n\beta)$

**Reduce-scatter**

$2log(p)\alpha + log(p)n(\beta + \gamma) + \dfrac{p-1}{p}n\beta$

**Allreduce**

**Allgather**

236

# Recap

**Reduce(-to-one)**

$log(p)(\alpha + n\beta + n\gamma)$

**Scatter**

$log(p)\alpha + \dfrac{p-1}{p}n\beta$

**Gather**

$log(p)\alpha + \dfrac{p-1}{p}n\beta$

**Broadcast**

$log(p)(\alpha + n\beta)$

**Reduce-scatter**

$2log(p)\alpha + log(p)n(\beta + \gamma) + \dfrac{p-1}{p}n\beta$

**Allreduce**

$2log(p)\alpha + log(p)n(2\beta + \gamma)$

**Allgather**

$2log(p)\alpha + log(p)n\beta + \dfrac{p-1}{p}n\beta$

237

# Recap

**Reduce(-to-one)**

$log(p)(\alpha + n\beta + n\gamma)$

**Scatter**

$log(p)\alpha + \dfrac{p-1}{p}n\beta$

**Gather**

$log(p)\alpha + \dfrac{p-1}{p}n\beta$

**Broadcast**

$log(p)(\alpha + n\beta)$

**Reduce-scatter**

$2log(p)\alpha + log(p)n(\beta + \gamma) + \dfrac{p-1}{p}n\beta$

**Allreduce**

$2log(p)\alpha + log(p)n(2\beta + \gamma)$

**Allgather**

238

# Recap

**Reduce(-to-one)**

$log(p)(\alpha + n\beta + n\gamma)$

**Scatter**

$log(p)\alpha + \dfrac{p-1}{p}n\beta$

**Gather**

$log(p)\alpha + \dfrac{p-1}{p}n\beta$

**Broadcast**

$log(p)(\alpha + n\beta)$

**Reduce-scatter**

$2log(p)\alpha + log(p)n(\beta + \gamma) + \dfrac{p-1}{p}n\beta$

**Allreduce**

$2log(p)\alpha + log(p)n(2\beta + \gamma)$

**Allgather**

$2log(p)\alpha + log(p)n\beta + \dfrac{p-1}{p}n\beta$

239

# A building block approach to library implementation

- **Short vector case**

- **Long vector case**

- **Hybrid algorithms**

# Long vector case

- **Primary concern:**
  - algorithms must have low cost due to vector length
  - algorithms must avoid network conflicts

- **Secondary concerns:**
  - algorithms must work for arbitrary number of nodes
    - » in particular, not just for power-of-two numbers of nodes

# Long vector building blocks

- **We will show how the following building blocks:**
  - collect/distributed combine
  - scatter/gather

  **can be implemented using "bucket" algorithms while attaining**
  - minimal cost due to length of vectors
  - implementation for arbitrary numbers of nodes
  - no network conflicts

- **NOTICE: scatter and gather already satisfy these conditions**

242

# General principles

- **A logical ring can be embedded in a physical linear array with worm-hole routing, since the "wrap-around" message doesn't conflict**
  - **This is used to "drop off" messages or to "pick up" contributions**

- **A logical ring can be embedded in a physical linear array with worm-hole routing, since the "wrap-around" message doesn't conflict**

245

247

# General principles

- Can be used to implement the following building blocks:
  - collect
  - distributed combine

  using a bucket algorithm embedded in the physical linear array while attaining
  - minimal cost due to vector length
  - implementation for arbitrary numbers of nodes
  - no network conflicts

# Allgather

**Before**

**After**



250

251

252

254

255

256

258

260

263

265

# Cost of bucket Allgather

$$(p-1)\left(\alpha + \frac{n}{p}\beta\right)$$

number of steps

cost per steps

$$=$$

$$(p-1)\alpha + \frac{p-1}{p}n\beta$$

269

# Cost of bucket Allgather

$$(p-1)\left(\alpha + \frac{n}{p}\beta\right)$$

number of steps

cost per steps

$$= (p-1)\alpha + \frac{p-1}{p}n\beta$$

Notice: attains lower bound for bandwidth component

# Reduce-scatter

**Before**

**After**



271

273

274

281

287

# Cost of bucket distributed combine

$$(p-1)\left(\alpha + \frac{n}{p}\beta + \frac{n}{p}\beta\right)$$

number of steps

cost per steps

$$= (p-1)\alpha + \frac{p-1}{p}n\beta + \frac{p-1}{p}n\gamma$$

290

# Cost of bucket Reduce-scatter

$$(p-1) \left( \alpha + \frac{n}{p}\beta + \frac{n}{p}\gamma \right)$$

$$=$$

$$(p-1)\alpha + \frac{p-1}{p}n\beta + \frac{p-1}{p}n\gamma$$

number of steps

cost per steps

Notice: attains lower bound for bandwidth and computation component

291

# Scatter

**Notice: Scatter as implemented before was optimal in latency and bandwidth components**

**Before**

**After**

# Gather

**Notice: Gather as implemented before was optimal in latency and bandwidth components**

**Before**

**After**

293

# Using the building blocks

# Broadcast (long vector)

# Broadcast (long vector)



Scatter

296

# Broadcast (long vector)



Allgather

297

# Cost of scatter/allgather broadcast

- **Assumption: power of two number of nodes**

**scatter**

**allgather**

$$\frac{\begin{array}{c} log(p)\alpha + \dfrac{p-1}{p}n\beta \\[2em] (p-1)\alpha + \dfrac{p-1}{p}n\beta \end{array}}{(log(p)+p-1)\alpha + 2\dfrac{p-1}{p}n\beta}$$

298

# Cost of scatter/allgather broadcast

- **Assumption: power of two number of nodes**

**scatter**

**allgather**

$$\frac{\begin{array}{l} log(p)\alpha + \dfrac{p-1}{p}n\beta \\[4mm] (p-1)\alpha + \dfrac{p-1}{p}n\beta \end{array}}{(log(p)+p-1)\alpha + 2\dfrac{p-1}{p}n\beta}$$

Notice: attains within a factor of two of the lower bound for bandwidth

299

# Reduce(-to-one) (long vector)

# Combine-to-one (long vector)



Reduce-scatter

301

# Combine-to-one (long vector)



Gather

302

# Cost of Reduce-scatter/Gather Reduce(-to-one)

- **Assumption: power of two number of nodes**

**Reduce-scatter**
$$(p-1)\alpha + \frac{p-1}{p}n\beta + \frac{p-1}{p}n\gamma$$

**gather**
$$\frac{log(p)\alpha + \frac{p-1}{p}n\beta}{(log(p)+p-1)\alpha + 2\frac{p-1}{p}n\beta + \frac{p-1}{p}n\gamma}$$

# Cost of Reduce-scatter/Gather Reduce(-to-one)

- **Assumption: power of two number of nodes**

**Reduce-scatter**
$$(p-1)\alpha + \frac{p-1}{p}n\beta + \frac{p-1}{p}n\gamma$$

**gather**
$$log(p)\alpha + \frac{p-1}{p}n\beta$$

$$(log(p) + p - 1)\alpha + 2\frac{p-1}{p}n\beta + \frac{p-1}{p}n\gamma$$

Notice: attains **within a factor of two** of the lower bound for bandwidth and attains lower bound for computation

304

# Allreduce
# (long vector)

# Allreduce
# (long vector)



**Reduce-scatter**

306

# Allreduce
# (long vector)



**Allgather**

307

# Cost of Reduce-scatter/Allgather Allreduce

- **Assumption: power of two number of nodes**

**Reduce-scatter** $(p-1)\alpha + \dfrac{p-1}{p}n\beta + \dfrac{p-1}{p}n\gamma$

**Allgather** $\dfrac{(p-1)\alpha + \dfrac{p-1}{p}n\beta}{2(p-1)\alpha + 2\dfrac{p-1}{p}n\beta + \dfrac{p-1}{p}n\gamma}$

308

# Cost of Reduce-scatter/Allgather Allreduce

- **Assumption: power of two number of nodes**

**Reduce-scatter**
$$(p-1)\alpha + \frac{p-1}{p}n\beta + \frac{p-1}{p}n\gamma$$

**Allgather**
$$\frac{(p-1)\alpha + \frac{p-1}{p}n\beta}{2(p-1)\alpha + 2\frac{p-1}{p}n\beta + \frac{p-1}{p}n\gamma}$$

**Notice: attains the lower bound for bandwidth and computation**

309

# Recap

**Reduce-scatter**
$$(p-1)\alpha + \frac{p-1}{p}n(\beta+\gamma)$$

**Scatter**
$$log(p)\alpha + \frac{p-1}{p}n\beta$$

**Gather**
$$log(p)\alpha + \frac{p-1}{p}n\beta$$

**Allgather**
$$(p-1)\alpha + \frac{p-1}{p}n\beta$$

**Reduce(-to-one)**

**Allreduce**

**Broadcast**

310

# Recap

**Reduce-scatter**
$$(p-1)\alpha + \frac{p-1}{p}n(\beta+\gamma)$$

**Scatter**
$$log(p)\alpha + \frac{p-1}{p}n\beta$$

**Gather**
$$log(p)\alpha + \frac{p-1}{p}n\beta$$

**Allgather**
$$(p-1)\alpha + \frac{p-1}{p}n\beta$$

**Reduce(-to-one)**
$$(p-1+log(p))\alpha + \frac{p-1}{p}n(2\beta+\gamma)$$

**Allreduce**

**Broadcast**

# Recap

**Reduce-scatter**
$$(p-1)\alpha + \frac{p-1}{p}n(\beta+\gamma)$$

**Scatter**
$$log(p)\alpha + \frac{p-1}{p}n\beta$$

**Gather**
$$log(p)\alpha + \frac{p-1}{p}n\beta$$

**Allgather**
$$(p-1)\alpha + \frac{p-1}{p}n\beta$$

**Reduce(-to-one)**
$$(p-1+log(p))\alpha + \frac{p-1}{p}n(2\beta+\gamma)$$

**Allreduce**
$$2(p-1)\alpha + \frac{p-1}{p}n(2\beta+\gamma)$$

**Broadcast**
$$(log(p)+p-1)\alpha + 2\frac{p-1}{p}n\beta$$

# Recap

**Reduce-scatter**
$$(p-1)\alpha + \frac{p-1}{p}n(\beta+\gamma)$$

**Scatter**
$$log(p)\alpha + \frac{p-1}{p}n\beta$$

**Gather**
$$log(p)\alpha + \frac{p-1}{p}n\beta$$

**Allgather**
$$(p-1)\alpha + \frac{p-1}{p}n\beta$$

**Reduce(-to-one)**
$$(p-1+log(p))\alpha + \frac{p-1}{p}n(2\beta+\gamma)$$

**Allreduce**
$$2(p-1)\alpha + \frac{p-1}{p}n(2\beta+\gamma)$$

**Broadcast**

313

# Recap

**Reduce-scatter**

$$(p-1)\alpha + \frac{p-1}{p}n(\beta+\gamma)$$

**Scatter**

$$log(p)\alpha + \frac{p-1}{p}n\beta$$

**Gather**

$$log(p)\alpha + \frac{p-1}{p}n\beta$$

**Allgather**

$$(p-1)\alpha + \frac{p-1}{p}n\beta$$

**Reduce(-to-one)**

$$(p-1+log(p))\alpha + \frac{p-1}{p}n(2\beta+\gamma)$$

**Allreduce**

$$2(p-1)\alpha + \frac{p-1}{p}n(2\beta+\gamma)$$

**Broadcast**

$$(log(p)+p-1)\alpha + 2\frac{p-1}{p}n\beta$$

# Advanced Techniques:

# Taking advantage of higher dimensions

315

# Physical 2D meshes

- **Simple solution: embed logical linear array**
  - problem: large $p$ implies high latency for bucket algorithms

- **Advanced solution: perform operation in each dimension**
  - collect:

    collect within rows, followed by collect within columns
  - distributed combine:

    same, in reverse

# Example: 2D Allgather

# Example: 2D Allgather



**Allgather in rows**

318

# Example: 2D Collect



**Allgather in columns**

319

# Cost of 2D Allgather

**row Allgather**

**column Allgather**

$$(c - 1)\alpha + (c - 1)\frac{n}{p}\beta$$

$$\frac{(r - 1)\alpha + (r - 1)\frac{c}{p}n\beta}{(r + c - 2)\alpha + \frac{p - 1}{p}n\beta}$$

**latency term is reduced**

**bandwidth term is unaffected**

320

# Example: 2D Scatter/Allgather Broadcast

# Example: 2D Scatter/Allgather Broadcast

**scatter in columns**

# Example: 2D Scatter/ Allgather Broadcast



Scatter in rows

# Example: 2D Scatter/ Allgather Broadcast

**Allgather in rows**

# Example: 2D Scatter/Collect Broadcast

Allgather
in columns

# Cost of 2D scatter/Allgather broadcast

$$(log(p) + r + c - 2)\alpha + 2\frac{p-1}{p}n\beta$$

# A building block approach to library implementation

- **Short vector case**

- **Long vector case**

- **Hybrid algorithms**

# Hybrid algorithms (intermediate length case)

- **algorithms must balance latency, cost due to vector length, and network conflicts**

# Example

- **We will illustrate the techniques using the broadcast as an example**
  - short vector:  minimum spanning tree broadcast

# Example: 2D Broadcast

# Example: 2D Broadcast



- **Option 1:**
  - MST broadcast in column
  - MST broadcast in rows

331

# Example: 2D Broadcast



- **Option 1:**
  - MST broadcast in column
  - MST broadcast in rows

# Example: 2D Broadcast



- **Option 2:**
  - Scatter in column
  - MST broadcast in rows
  - Allgather in columns

333

# Example: 2D Broadcast



- **Option 2:**
  - Scatter in column
  - **MST broadcast in rows**
  - Allgather in columns

# Example: 2D Broadcast



- **Option 2:**
  - Scatter in column
  - MST broadcast in rows
  - Allgather **in columns**

335

# Example: 2D Broadcast



- **Option 3:**
    - **Scatter in column**
    - Scatter in rows
    - Allgather in rows
    - Allgather in columns

336

# Example: 2D Broadcast



- **Option 3:**
  - Scatter in column
  - **Scatter in rows**
  - Allgather in rows
  - Allgather in columns

# Example: 2D Broadcast



- **Option 3:**
  - Scatter in column
  - Scatter in rows
  - Allgather **in rows**
  - Allgather in columns

338

# Example: 2D Broadcast



- **Option 3:**
  - **Scatter in column**
  - **Scatter in rows**
  - **Allgather in rows**
  - **Allgather in columns**

339

# Cost comparison

- **Option 1:**
  - **MST broadcast in column**
  - **MST broadcast in rows**

- **Option 2:**
  - **Scatter in column**
  - **MST broadcast in rows**
  - **Allgather in columns**

- **Option 3:**
  - **Scatter in column**
  - **Scatter in rows**
  - **Allgather in rows**
  - **Allgather in columns**

$$\frac{\begin{array}{c} log(c)\alpha + log(c)n\beta \\ log(r)\alpha + log(r)n\beta \end{array}}{log(p)\alpha + log(p)n\beta}$$

# Cost comparison

- **Option 1:**
  - MST broadcast in column
  - MST broadcast in rows
- **Option 2:**
  - Scatter in column
  - MST broadcast in rows
  - Allgather in columns
- **Option 3:**
  - Scatter in column
  - Scatter in rows
  - Allgather in rows
  - Allgather in columns

$$
log(c)\alpha + \frac{c-1}{c}n\beta
$$

$$
log(r)\alpha + log(r)\frac{n}{c}\beta
$$

$$
(c-1)\alpha + \frac{c-1}{c}n\beta
$$

$$
(log(p) + c - 1)\alpha + \left(2\frac{c-1+log(r)}{c}\right)n\beta
$$

# Cost comparison

- **Option 1:**
  - MST broadcast in column
  - MST broadcast in rows

- **Option 2:**
  - Scatter in column
  - MST broadcast in rows
  - Allgather in columns

- **Option 3:**
  - Scatter in column
  - Scatter in rows
  - Allgather in rows
  - Allgather in columns

$$\frac{\begin{array}{l} log(c)\alpha + \dfrac{c-1}{c}n\beta \\[1em] log(r)\alpha + \dfrac{r-1}{r}\dfrac{n}{c}\beta \\[1em] (r-1)\alpha + \dfrac{r-1}{r}\dfrac{n}{c}\beta \\[1em] (c-1)\alpha + \dfrac{c-1}{c}n\beta \end{array}}{(log(p)+r+c-2)\alpha + 2\dfrac{p-1}{p}n\beta}$$

342

# Cost comparison

- **Option 1:**
  - MST broadcast in column
  - MST broadcast in rows

$$log(p)\alpha + log(p)n\beta$$

- **Option 2:**
  - Scatter in column
  - MST broadcast in rows
  - Allgather in columns

$$\left(log(p) + c - 1\right)\alpha + \left(2\frac{c - 1 + log(r)}{c}\right)n\beta$$

- **Option 3:**
  - Scatter in column
  - Scatter in rows
  - Allgather in rows
  - Allgather in columns

$$(log(p) + r + c - 2)\alpha + 2\frac{p - 1}{p}n\beta$$

# Higher dimensions

- **This technique can be extended by viewing one- and two-dimensional meshes logically as higher dimensions**
  - reduces latency
  - incurs network conflicts
  - can be used to create faster short vector implementations

- **Details require more time that is available today**

344

# Other techniques

- **Pipelined algorithms**
  - can be used to further reduce the cost of broadcast and combine-to-one for long vectors
  - very effective on hypercubes
    - » (Ho and Johnsson)
  - effective on meshes with low latency
    - » (Watts and van de Geijn)
  - complicated to implement, analyze and explain

345

# Outline

**Part I:  Theory**

- **Model of computation**
- **Collective communications**
- **A building block approach to library implementation**

## Part II: Practice

- **Implementation on the Paragon**
- **Performance results**
- **Applications**

346

# Outline

**Part I: Theory**

- **Model of computation**

- **Collective communications**

- **A building block approach to library implementation**

**Part II: Practice**

- **Implementation on the Paragon**

- **Performance results**

- **Applications**

347

# Theory is nice, but how does it work in practice?

- **Paragon does not match our model**
  - **Bad news:**
    - » **sending and receiving more complex then the model indicates**
    - » **forced messages vs. unforced messages**
    - » **preposted messages vs. nonpreposted messages**
    - » **etc.**
  - **Good news:**
    - » **excess bandwidth in the network**

# Interprocessor Collective Communication (InterCom) Project

# Implementation on the Paragon

- **Short vector building blocks**
    - reduce latency by *not* preposting and synchronizing
- **Long vector building blocks**
    - improve bandwidth by preposting and synchronizing
- **Incorporate more complex issues into model**
    - various startups, bandwidths, depending on situation
- **Use simple heuristic to choose hybrid strategy**
    - because of excess bandwidth, the mesh acts more like a hypercube, for which some solid theory exists
        - » (van de Geijn)
    - details go beyond this tutorial.

# Performance

351

# Performance comparison

- **NX collective communication**
- **Message Passing Interface (MPI)**
  - **Reference implementation from ANL and MSU**
  - **Bill Gropp, Rusty Lusk, and Tony Skjellum**
- **Basic Linear Algebra Communication Subprograms (BLACS)**
  - **Communication library of ScaLAPACK**
  - **Reference implementation from the Univ. of TN**
  - **Jack Dongarra and Clint Whaley**
- **Interprocessor Collective Communication (iCC) Library**
  - **High performance implementation by the InterCom team**

**Broadcast on 16 x 32 mesh Paragon**

353

**Broadcast on 16 x 32 mesh Paragon**

354

# Broadcast on 16 x 32 mesh Paragon



Legend:
- BLACS
- iCC
- MPI
- NX

Y-axis: **Time (seconds)** — 1.0000, 0.1000, 0.0100, 0.0010, 0.0001

X-axis: **Message Length (bytes)** — 1.0E+00, 1.0E+02, 1.0E+04, 1.0E+06

BLACS

# Broadcast on 16 x 32 mesh Paragon



356

# Broadcast on 16 x 32 mesh Paragon



357

# Allgather on 16 x 32 mesh Paragon



358

# Allgather on 16 x 32 mesh Paragon



359
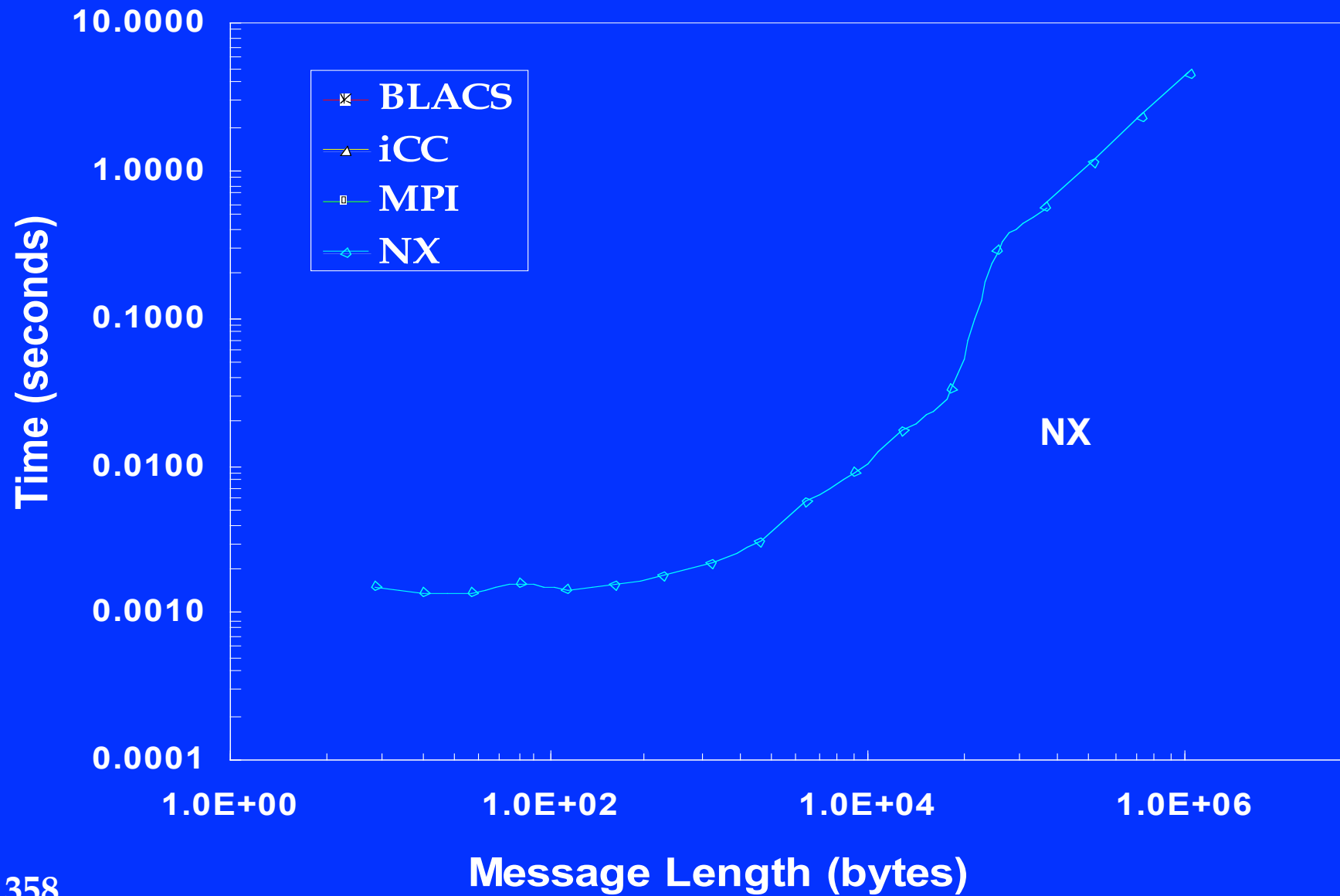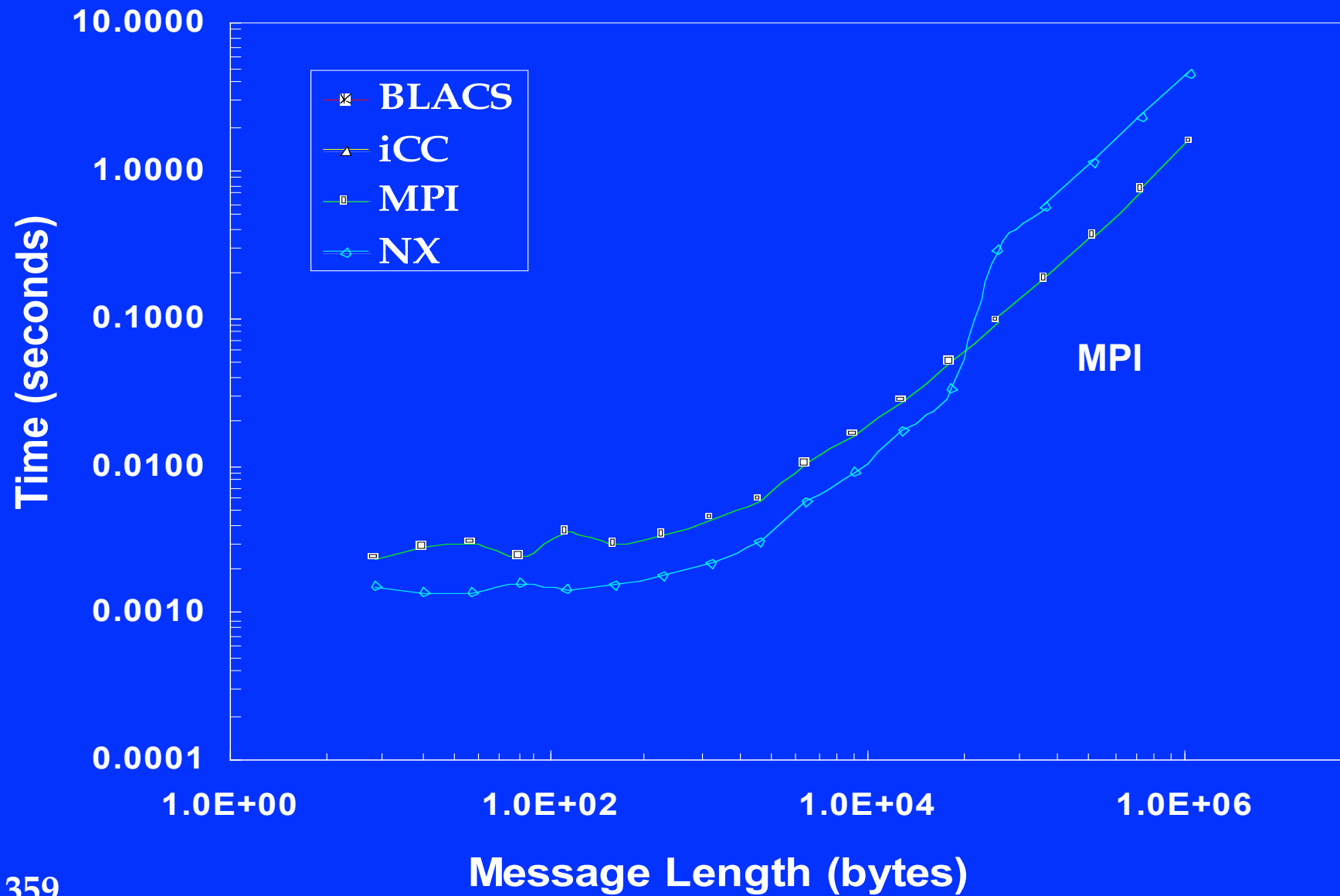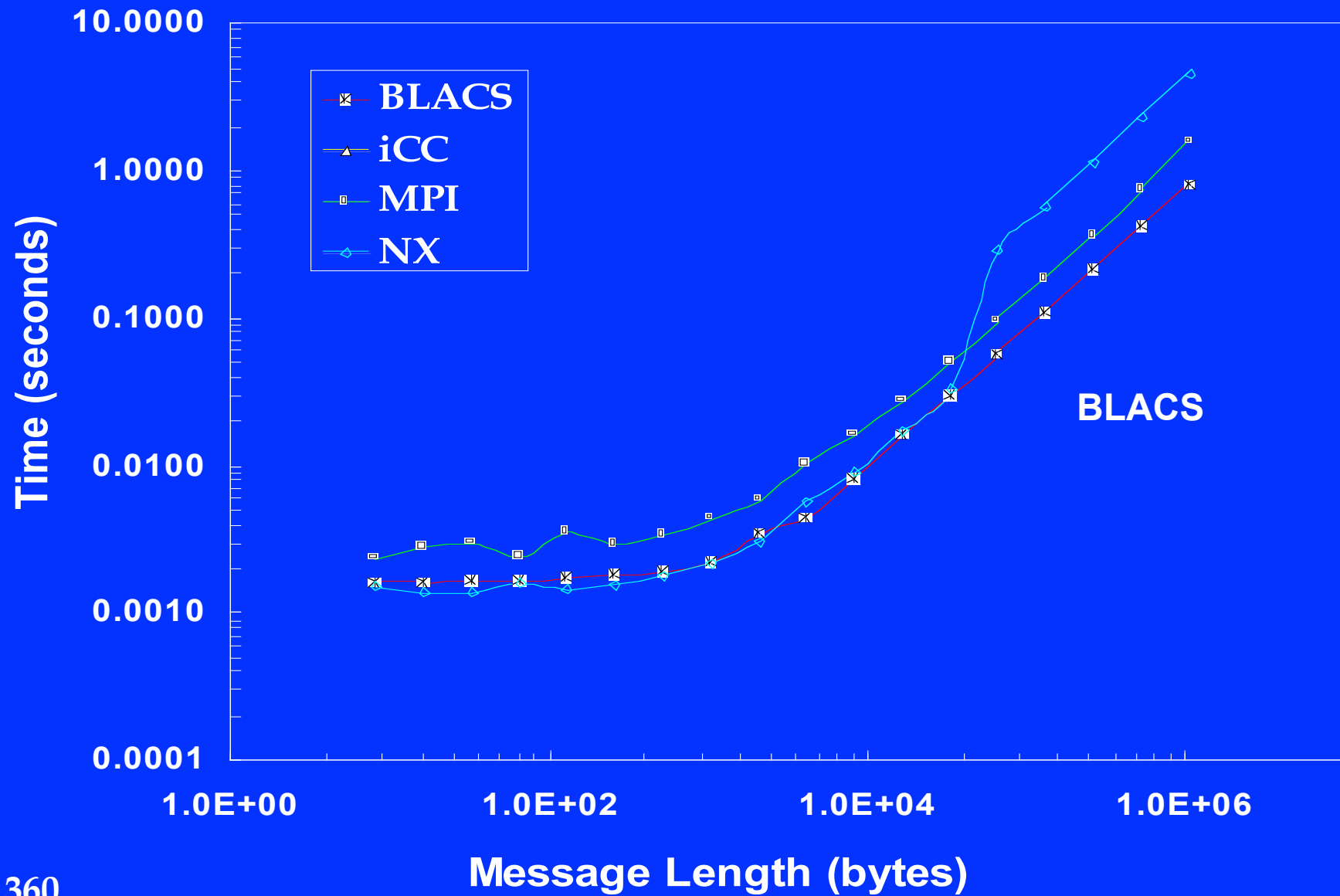
# Allgather on 16 x 32 mesh Paragon

# Allgather on 16 x 32 mesh Paragon



361

**Allgather on 16 x 32 mesh Paragon**

Legend:
- BLACS
- iCC
- MPI
- NX

Y-axis: Time (seconds), from 0.0 to 5.0

X-axis: Message Length (bytes), from 0.0E+00 to 1.2E+06

362

This PowerPoint presentation may be copied for nonprofit educational purposes.  Credit should be given to the InterCom project.

For information, contact

rvdg@cs.utexas.edu

363

# CollMark: Collective Communicaton Benchmark

## A look at the current state-of-the-art
## (spring 2000)

# How to measure the quality of an implementation

- **Architecture independent measure of the quality of the implementation:**

$$\frac{T_{comm}(n,p)}{T_{p2p}(n)}$$

- **Ideally:**

$$\frac{T_{comm}(n,p)}{T_{p2p}(n)} \xrightarrow[n\to\infty]{} 1 \qquad or \qquad 2$$

Broadcast 256 nodes