# Read-ME: Refactorizing LLMs as Router-Decoupled Mixture-of-Experts with System Co-design

**Friday, 13 Dec 11AM-2PM**
**Poster Session 5**

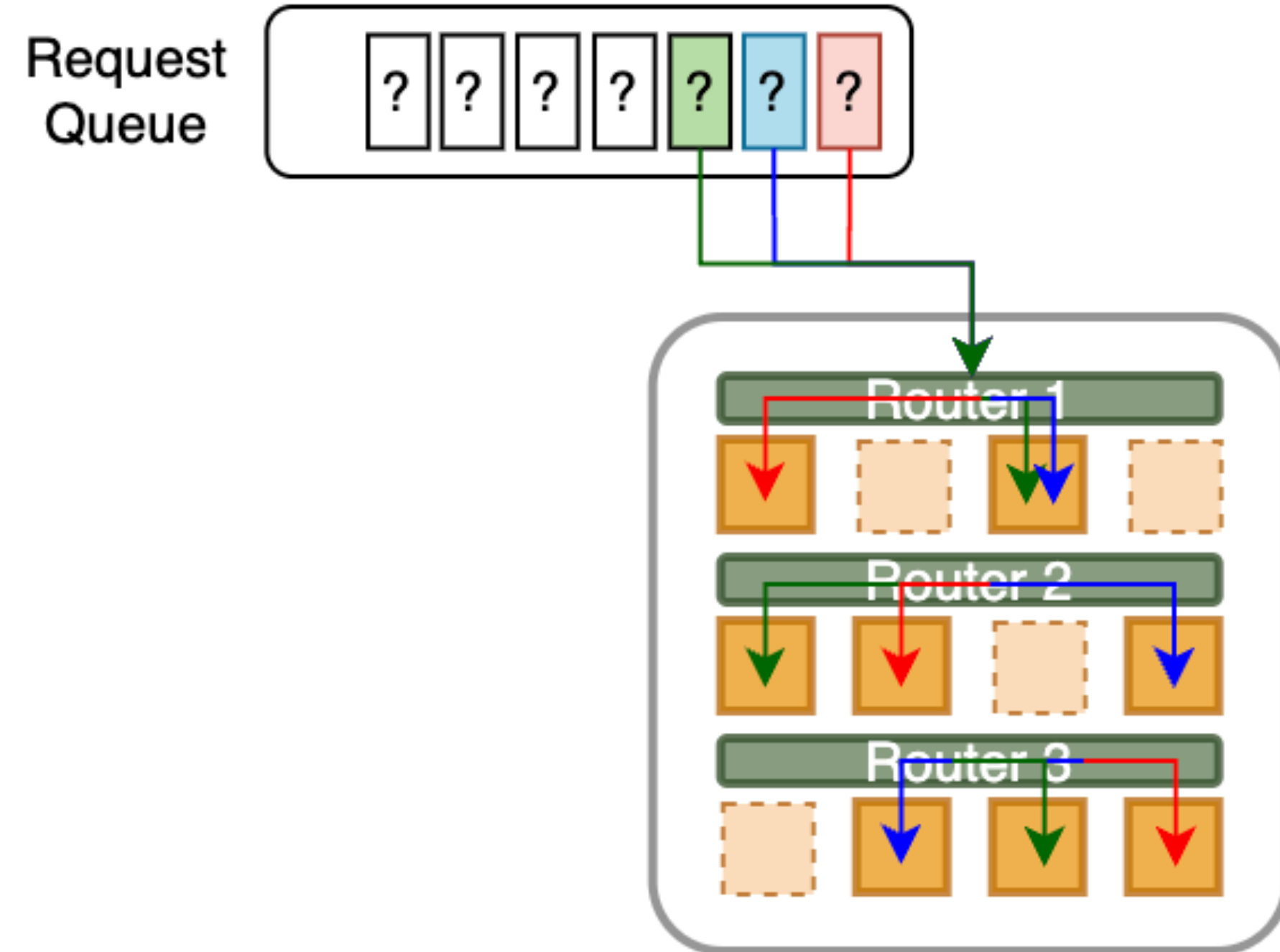Ruisi Cai*, **Yeonju Ro***, Geon-Woo Kim, Peihao Wang, Babak Ehteshami Bejnordi, Aditya Akella, Zhangyang Wang

Qualcomm   VITA   UTNS
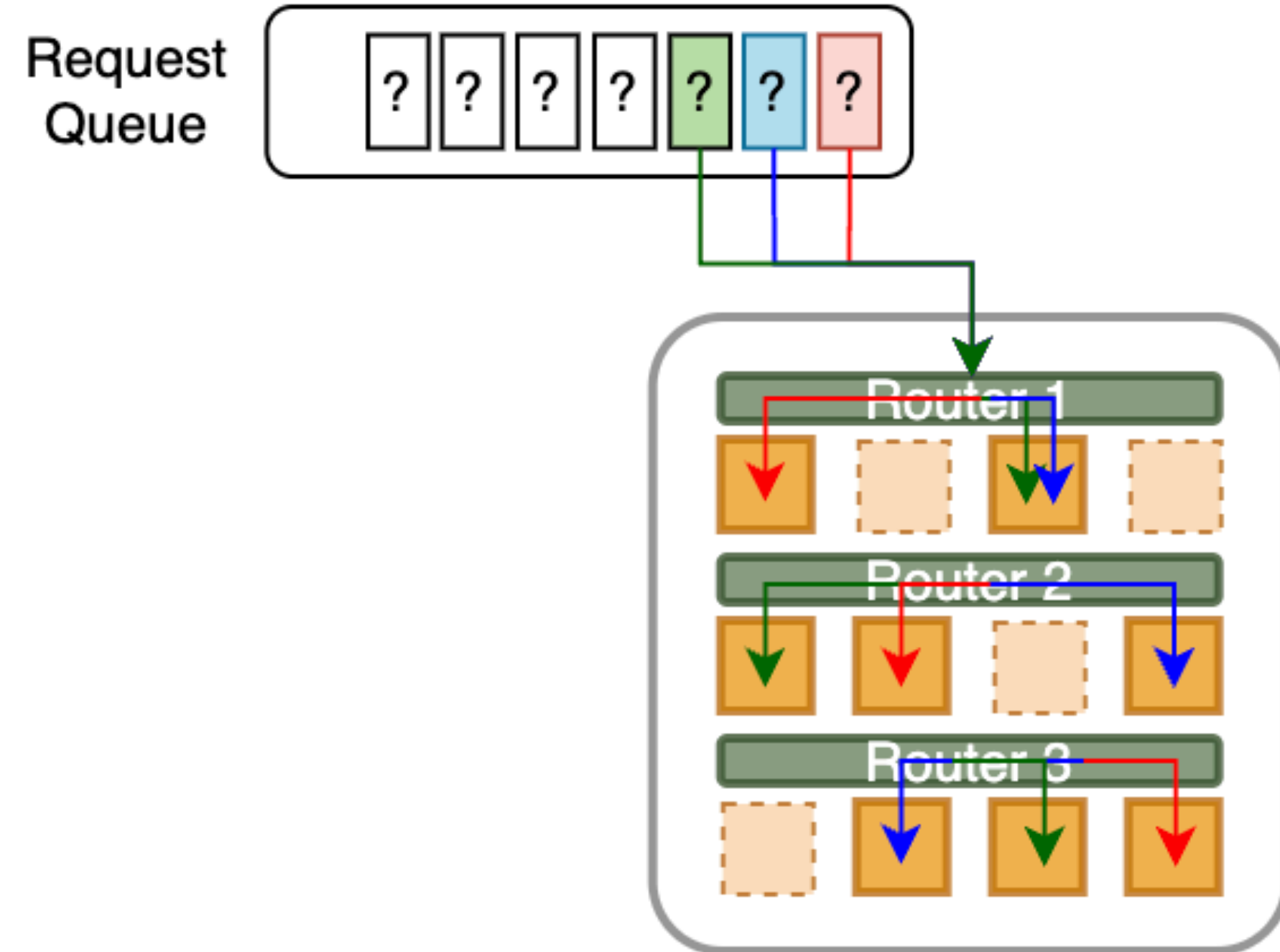
# Drawback of current MoE Architecture
## Layer-wise Gating Disrupts Efficient Memory Management
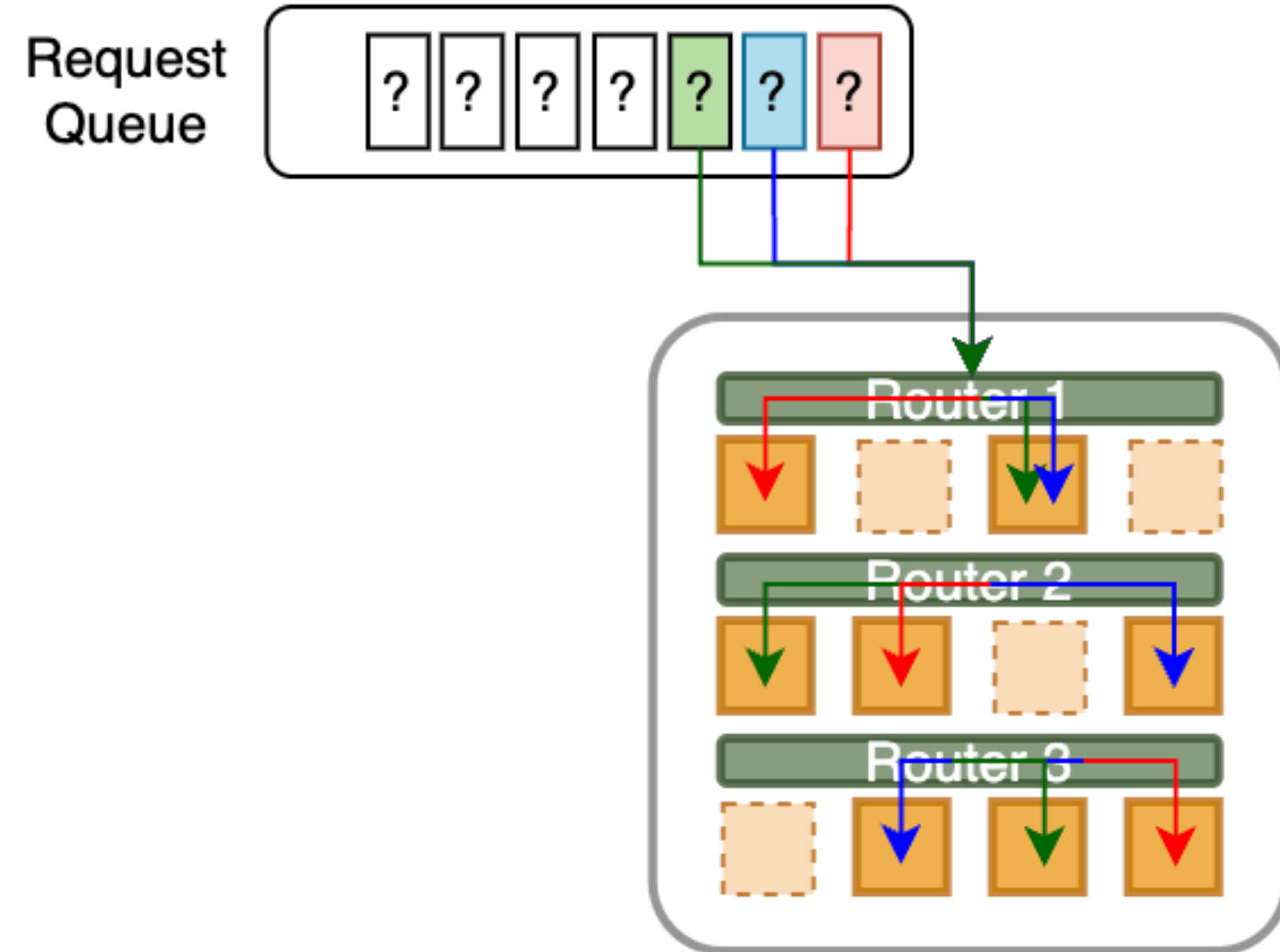
# Drawback of current MoE Architecture
## Layer-wise Gating Disrupts Efficient Memory Management



- Ideally, we only want to load layers to be used.
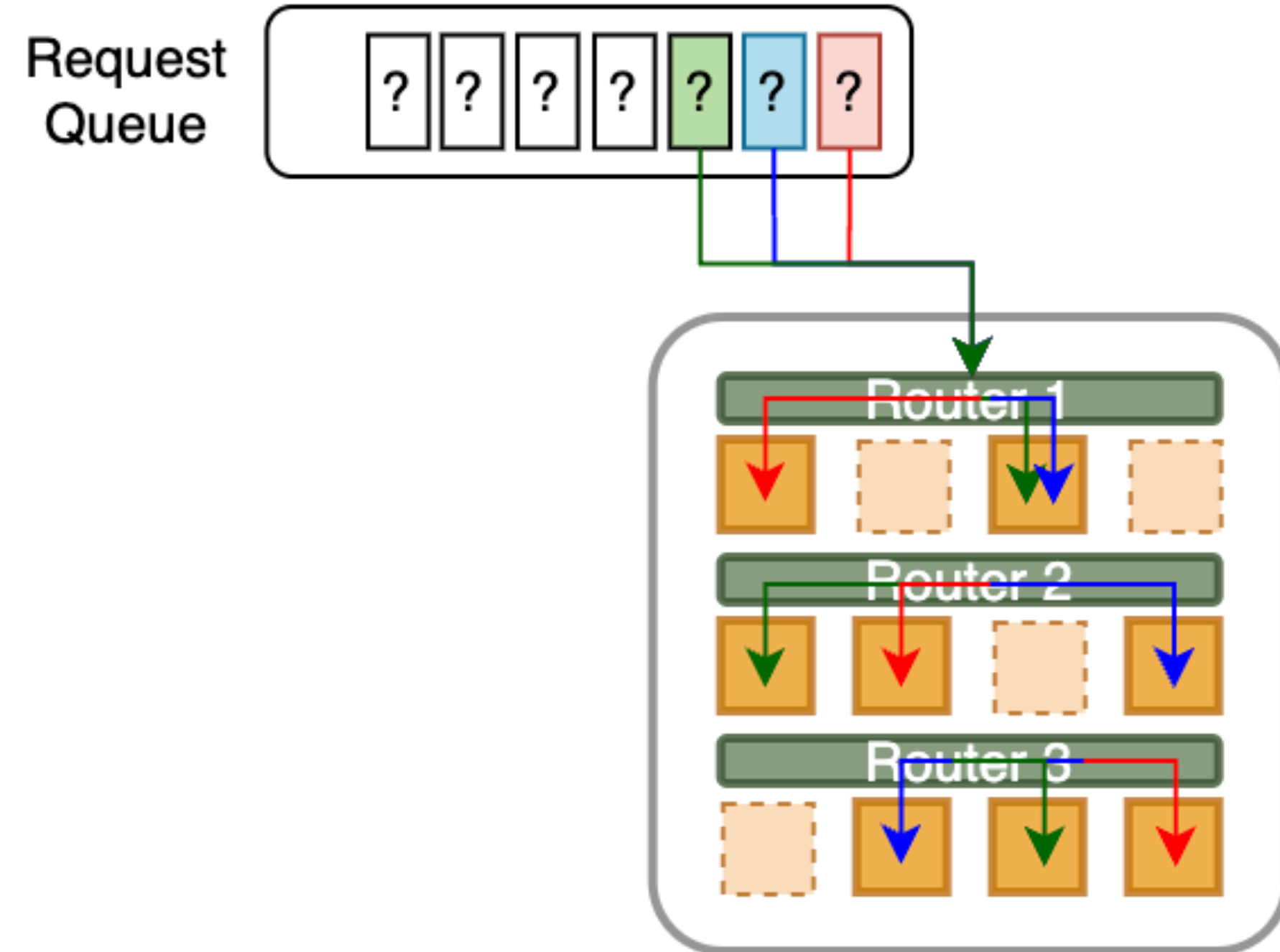
# Drawback of current MoE Architecture
## Layer-wise Gating Disrupts Efficient Memory Management



- Ideally, we only want to load layers to be used.

- However, we do not know which expert to activate until we go through the layer-wise routing layer.
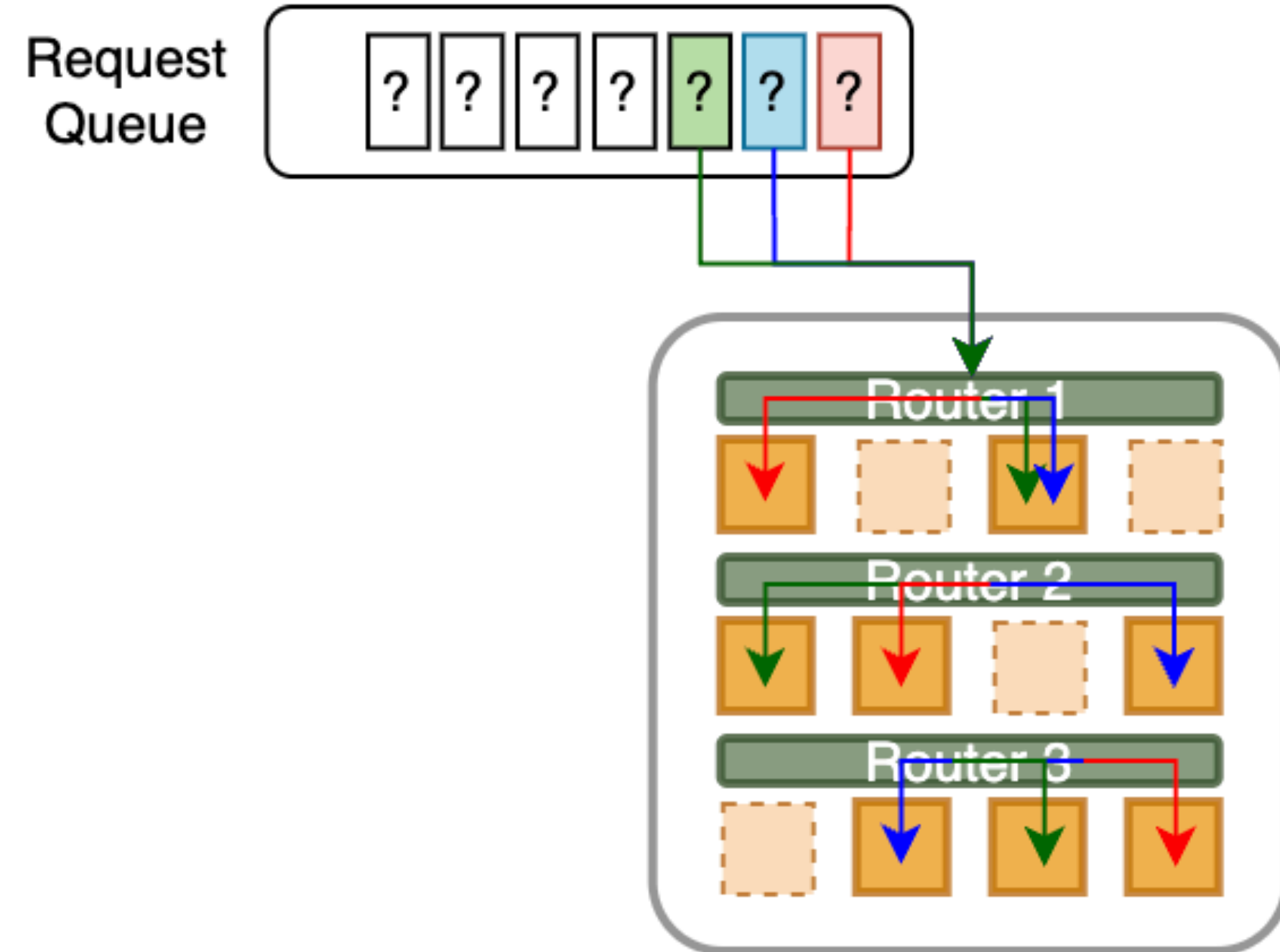
# Drawback of current MoE Architecture
## Layer-wise Gating Disrupts Efficient Memory Management



- Ideally, we only want to load layers to be used.

- However, we do not know which expert to activate until we go through the layer-wise routing layer.

- Thus, we need to either load all experts or load the expert to be used *on-demand*, which potentially add loading latency to the critical path.

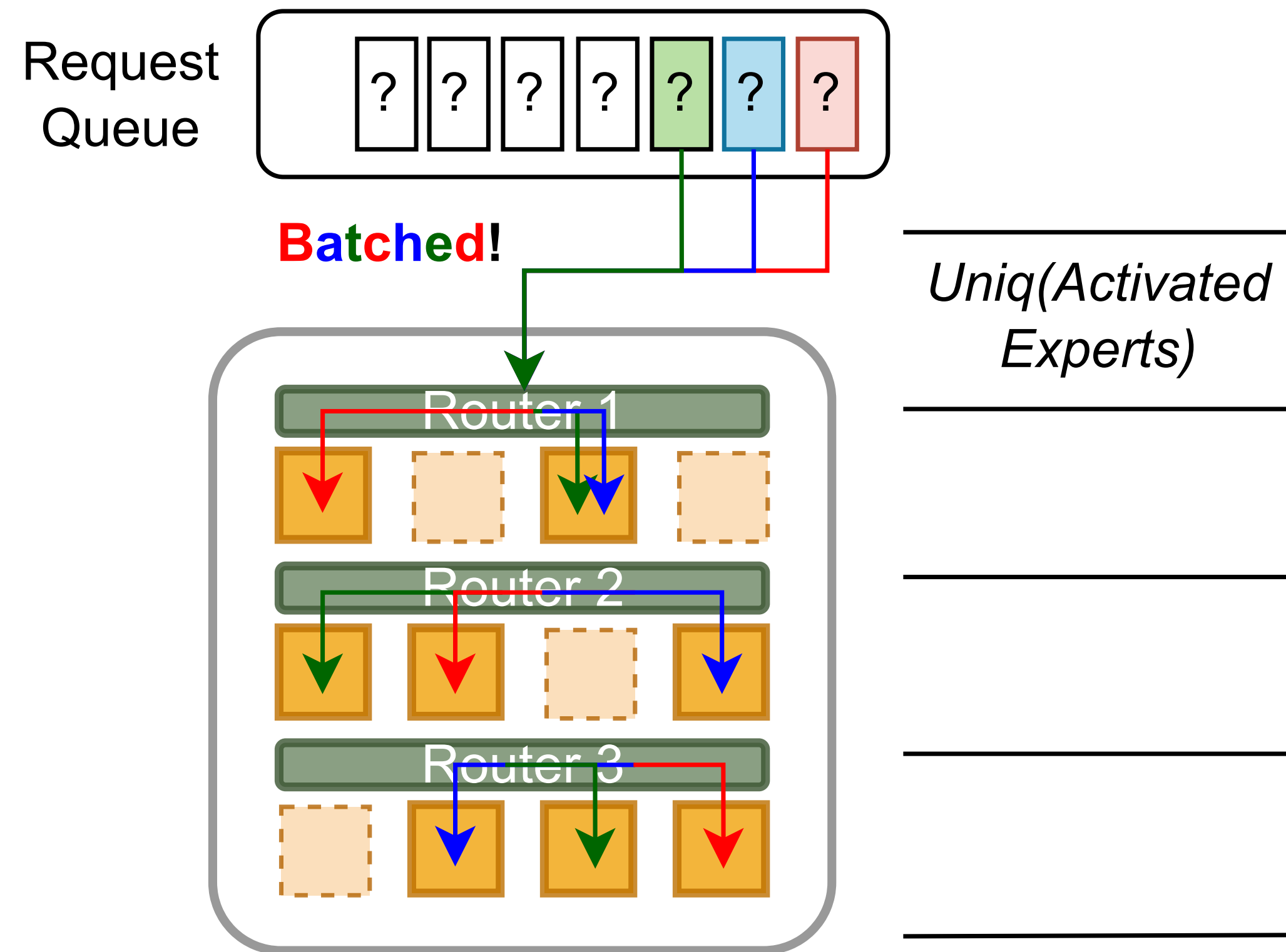# Drawback of current MoE Architecture
## Layer-wise Gating Disrupts Efficient Memory Management



- Ideally, we only want to load layers to be used.

- However, we do not know which expert to activate until we go through the layer-wise routing layer.

- Thus, we need to either load all experts or load the expert to be used *on-demand*, which potentially add loading latency to the critical path.

- Both prefetching/caching are not trivial.
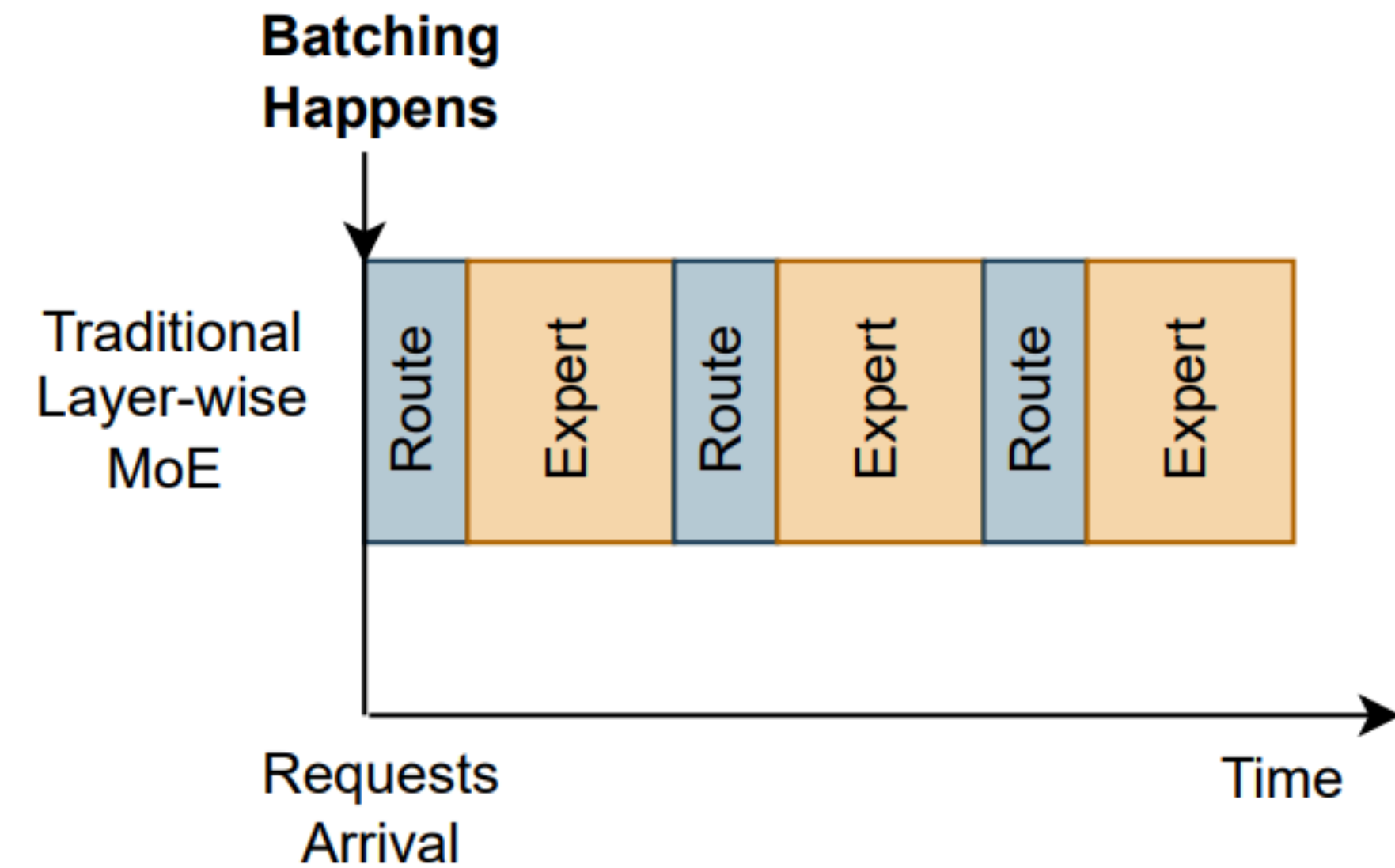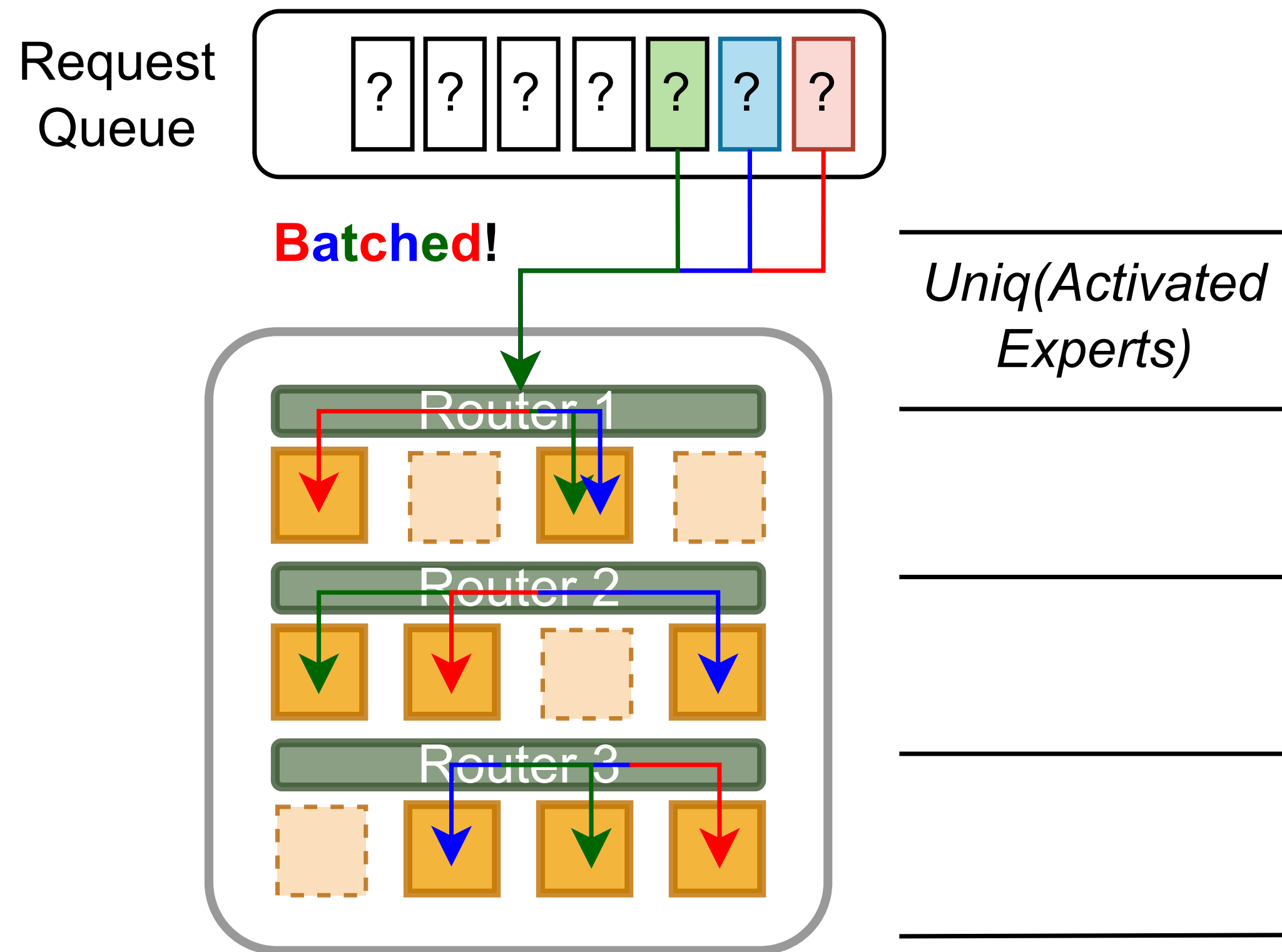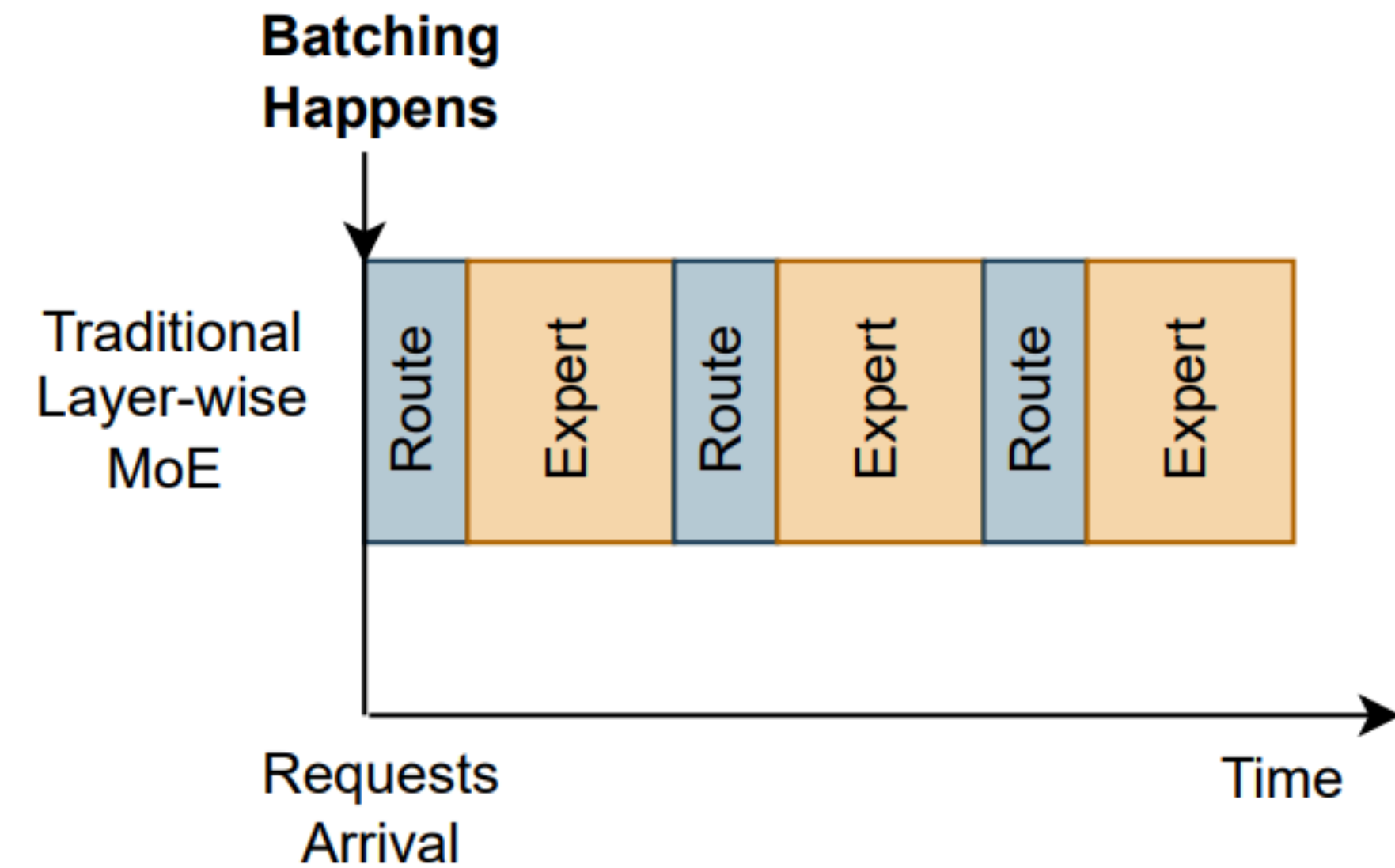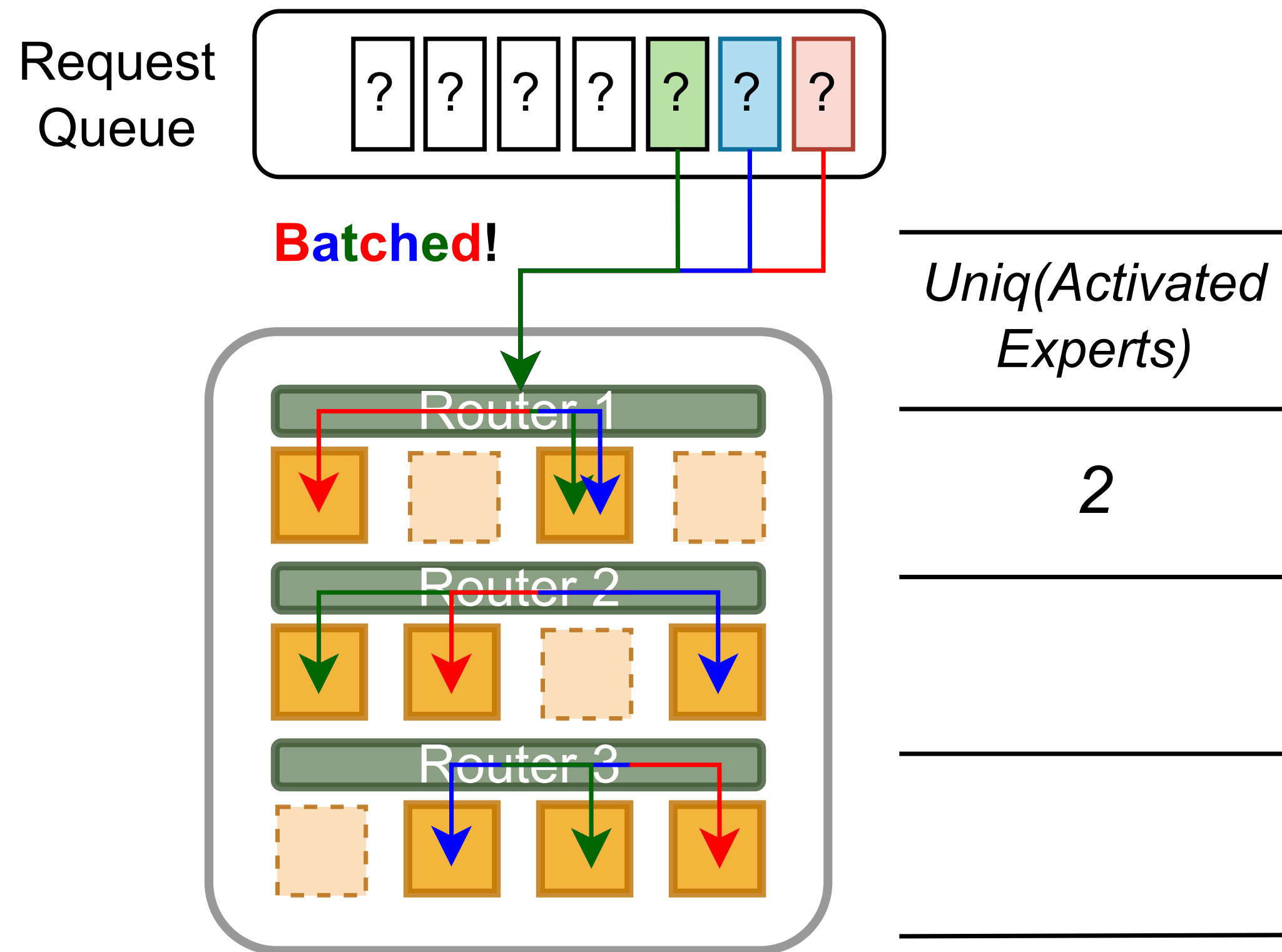
# Current Inference System
## Layer-wise Gating Disrupts Efficient Batching

# Current Inference System
## Layer-wise Gating Disrupts Efficient Batching
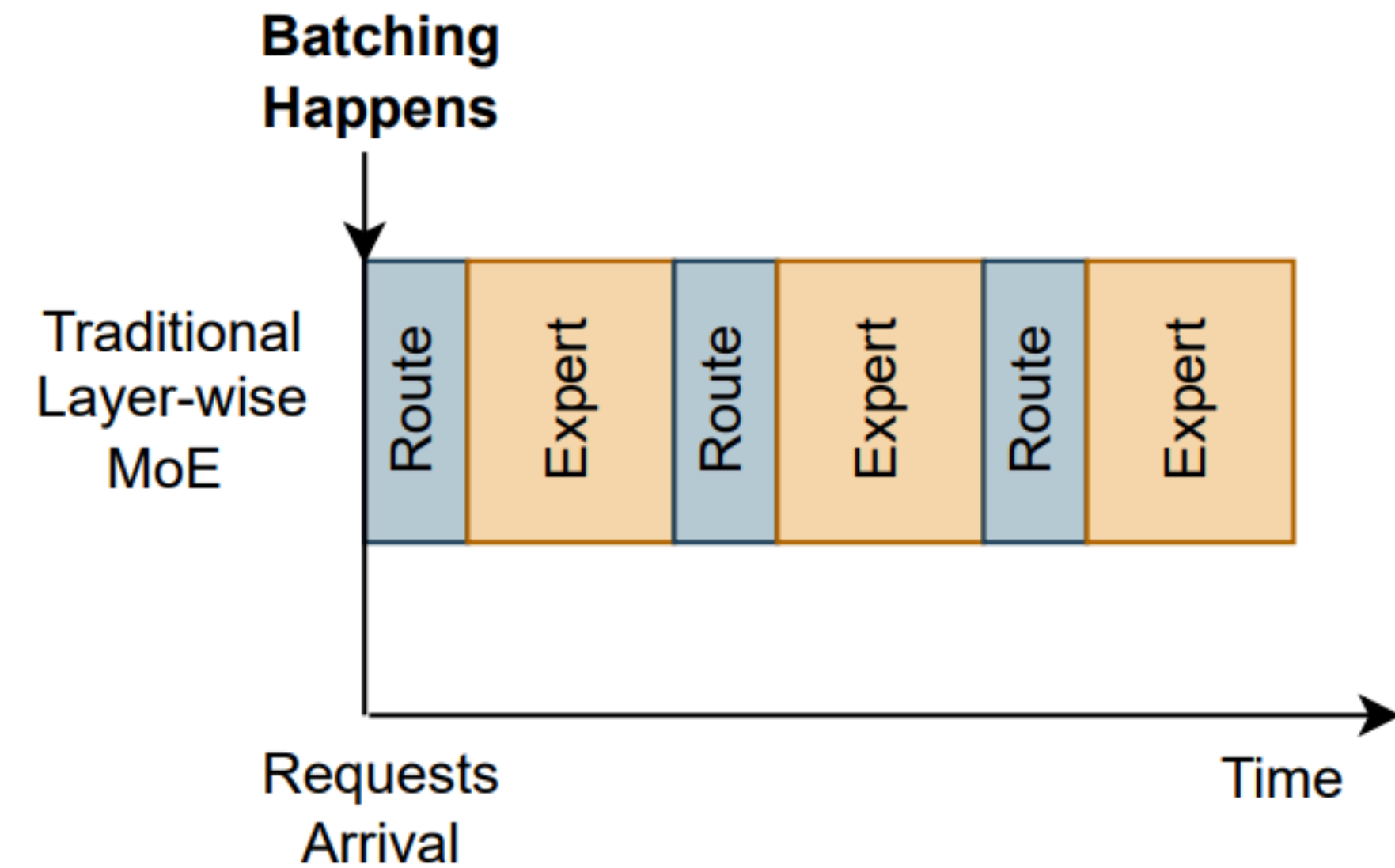
# Current Inference System
## Layer-wise Gating Disrupts Efficient Batching



Request Queue

? ? ? ? ? ? ?

Batched!

*Uniq(Activated Experts)*

Router 1

2

Router 2

Router 3

**Batching Happens**

Traditional Layer-wise MoE

Route | Expert | Route | Expert | Route | Expert
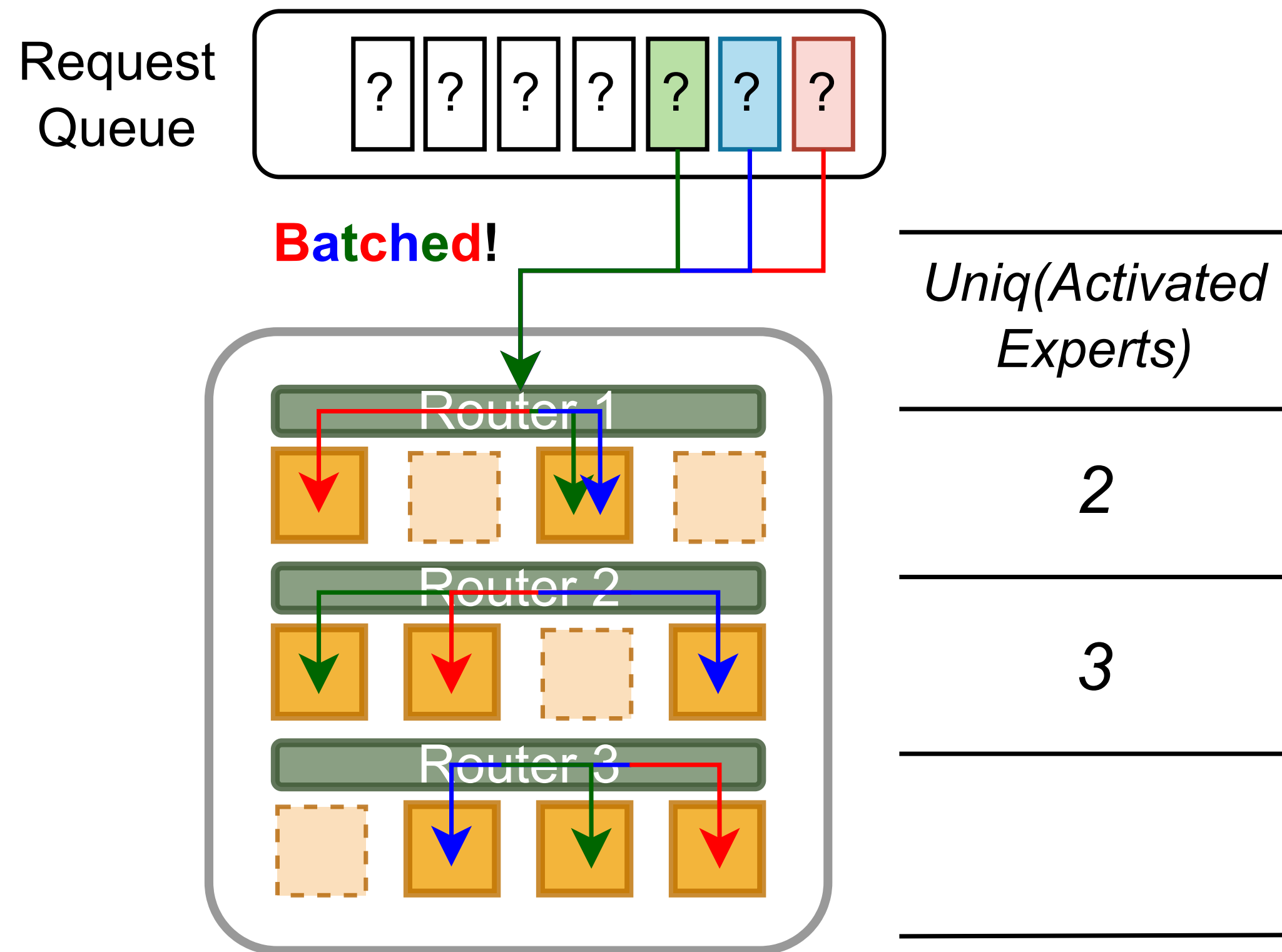
Requests Arrival

Time
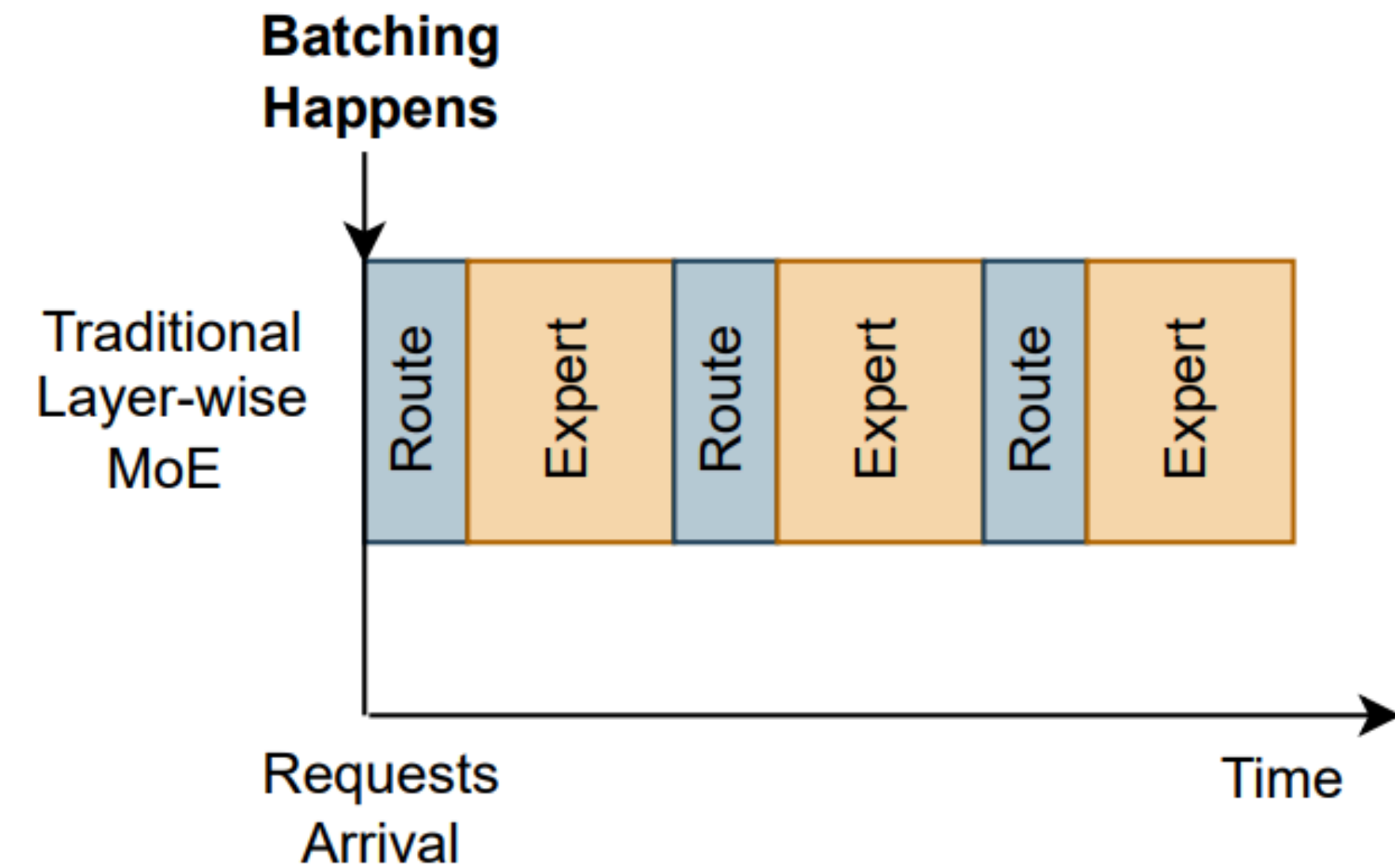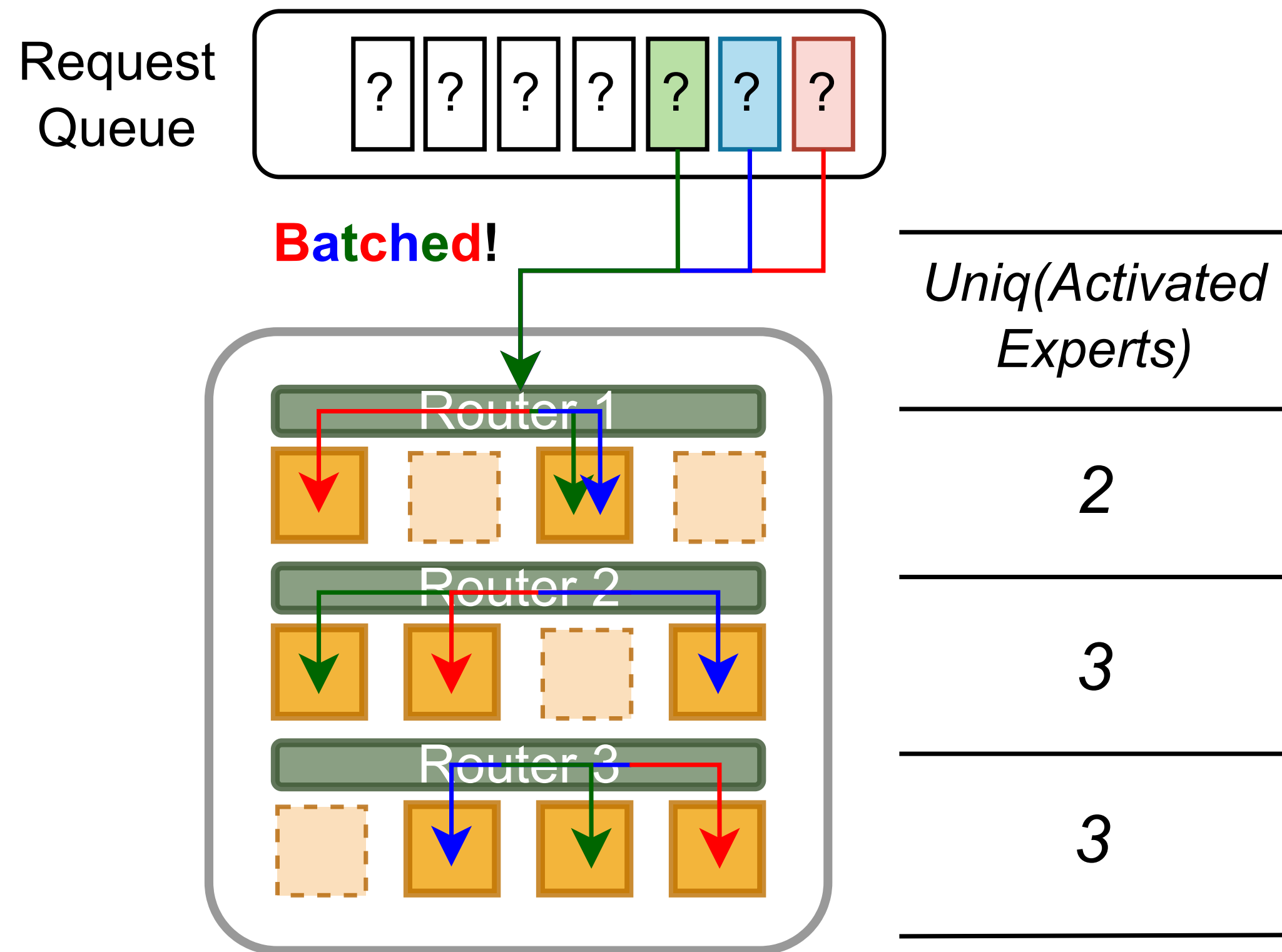
# Current Inference System
## Layer-wise Gating Disrupts Efficient Batching

# Current Inference System
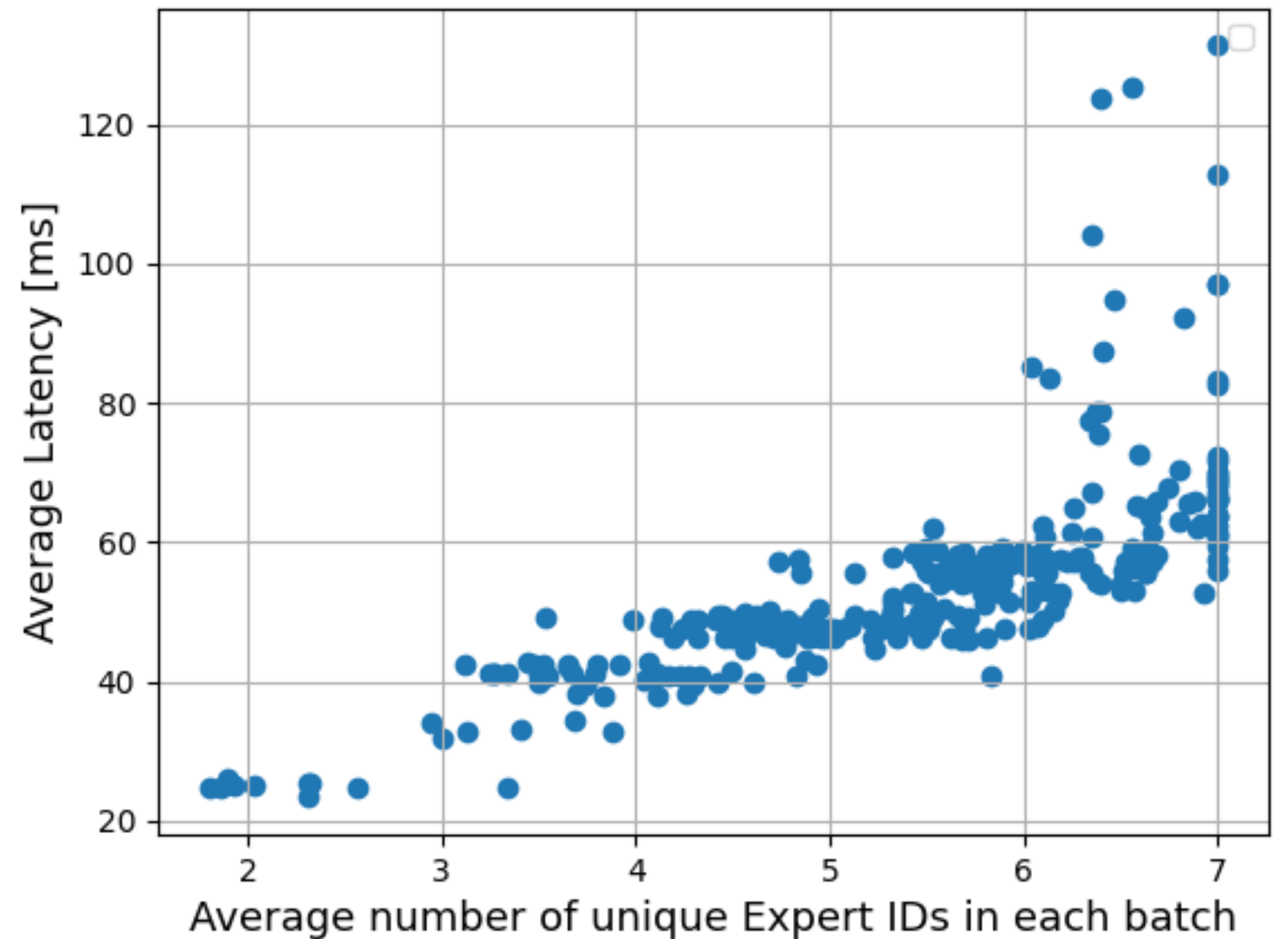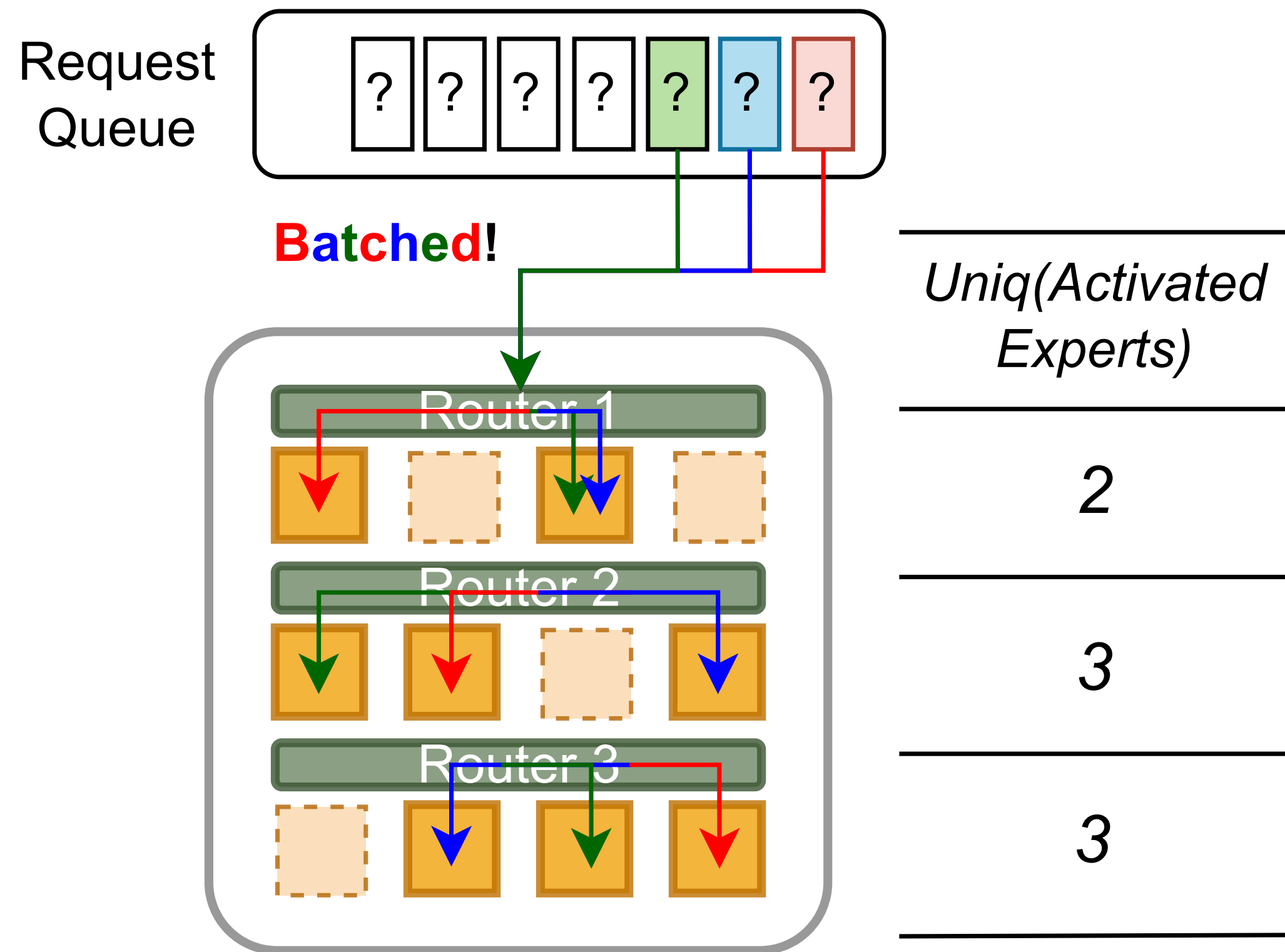## Layer-wise Gating Disrupts Efficient Batching



Request Queue

? ? ? ? ? ? ?

**Batched**!

Router 1

Router 2

Router 3

*Uniq(Activated Experts)*

2

3

3

**Batching Happens**

Traditional Layer-wise MoE

Route | Expert | Route | Expert | Route | Expert
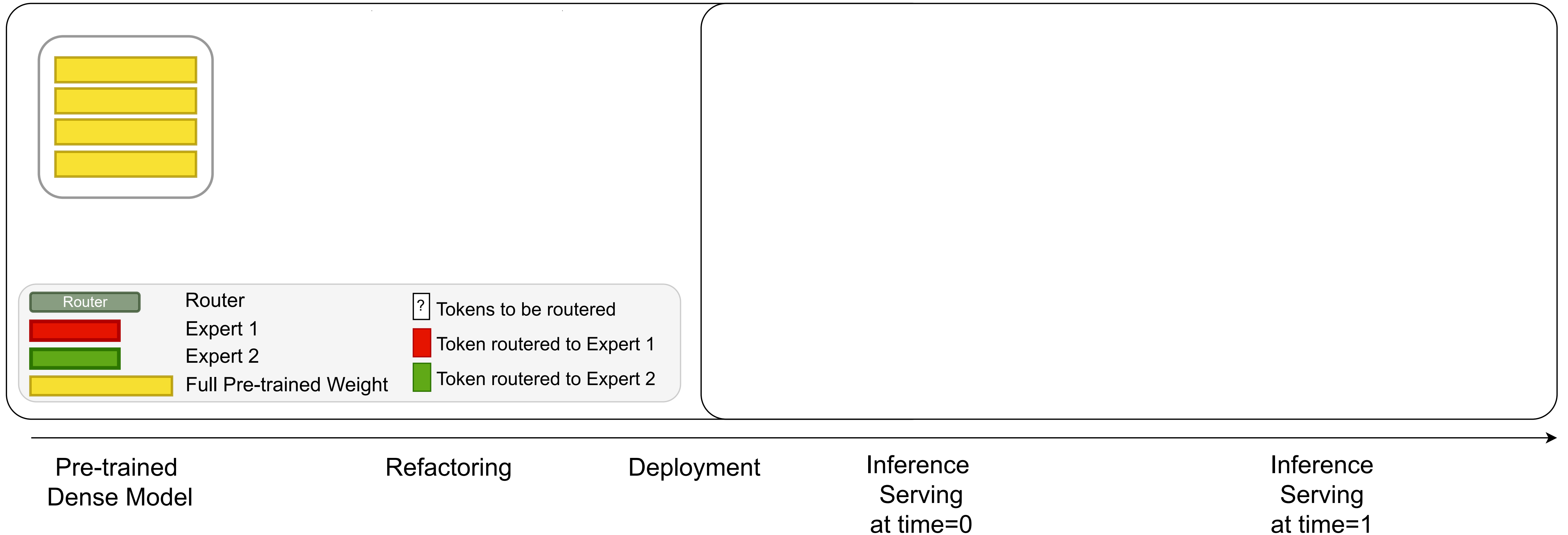
Requests Arrival

Time

# Current Inference System
## Token batching is poorly suited for MoE architecture

# Solution: Decoupled Router
## Re-designing MoE Architecture with Decoupled Router



| | | | |
|---|---|---|---|
| Router | Router | ? | Tokens to be routered |
| | Expert 1 | | Token routered to Expert 1 |
| | Expert 2 | | Token routered to Expert 2 |
| | Full Pre-trained Weight | | |

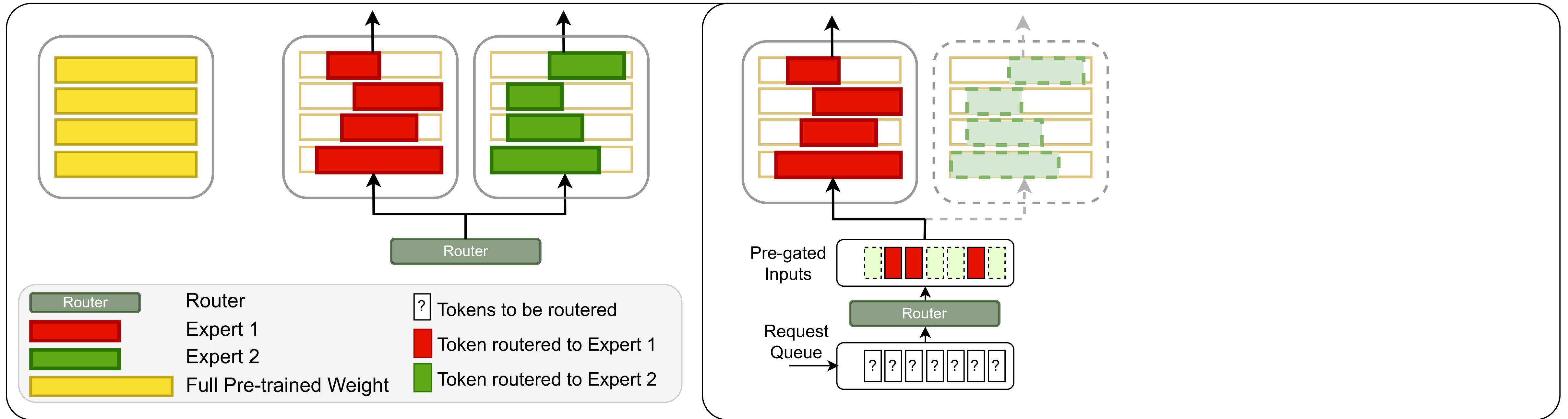Pre-trained
Dense Model

Refactoring

Deployment

Inference
Serving
at time=0

Inference
Serving
at time=1

# Solution: Decoupled Router
## Re-designing MoE Architecture with Decoupled Router



Router | Router
Expert 1
Expert 2
Full Pre-trained Weight

? | Tokens to be routered
Token routered to Expert 1
Token routered to Expert 2
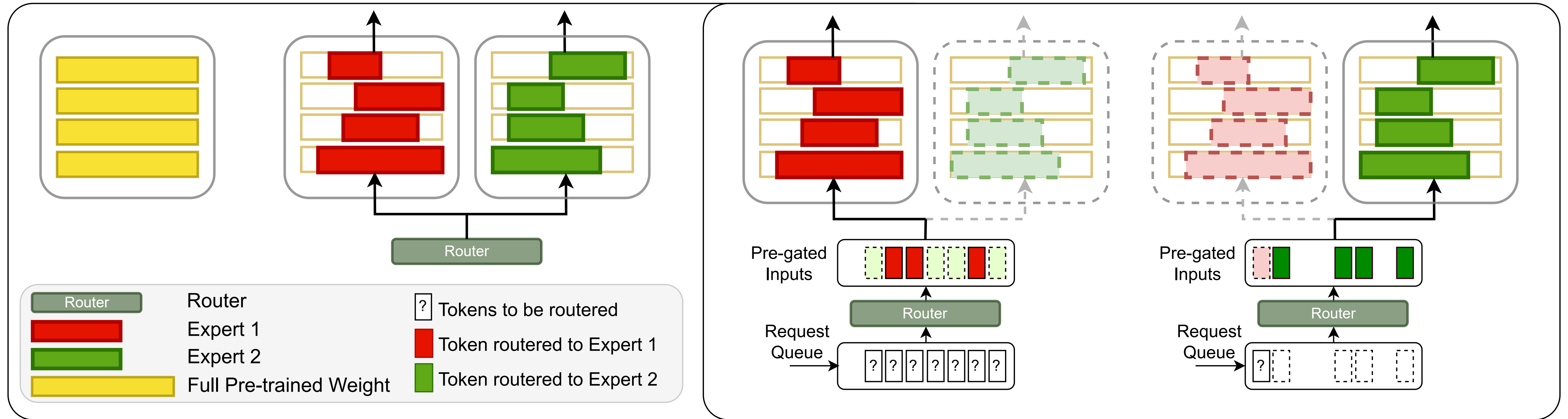
Pre-trained
Dense Model

Refactoring

Deployment

Inference
Serving
at time=0

Inference
Serving
at time=1

# Solution: Decoupled Router
## Re-designing MoE Architecture with Decoupled Router



| | | | | |
|---|---|---|---|---|
| Pre-trained<br>Dense Model | Refactoring | Deployment | Inference<br>Serving<br>at time=0 | Inference<br>Serving<br>at time=1 |

Legend:
- Router — Router
- Expert 1
- Expert 2
- Full Pre-trained Weight
- ? Tokens to be routered
- Token routered to Expert 1
- Token routered to Expert 2

# Solution: Decoupled Router
## Re-designing MoE Architecture with Decoupled Router



Router · Router
Expert 1
Expert 2
Full Pre-trained Weight

? Tokens to be routered
Token routered to Expert 1
Token routered to Expert 2

Pre-gated Inputs
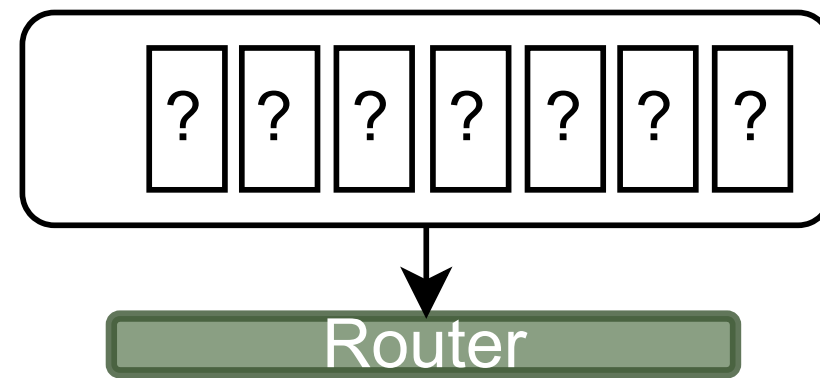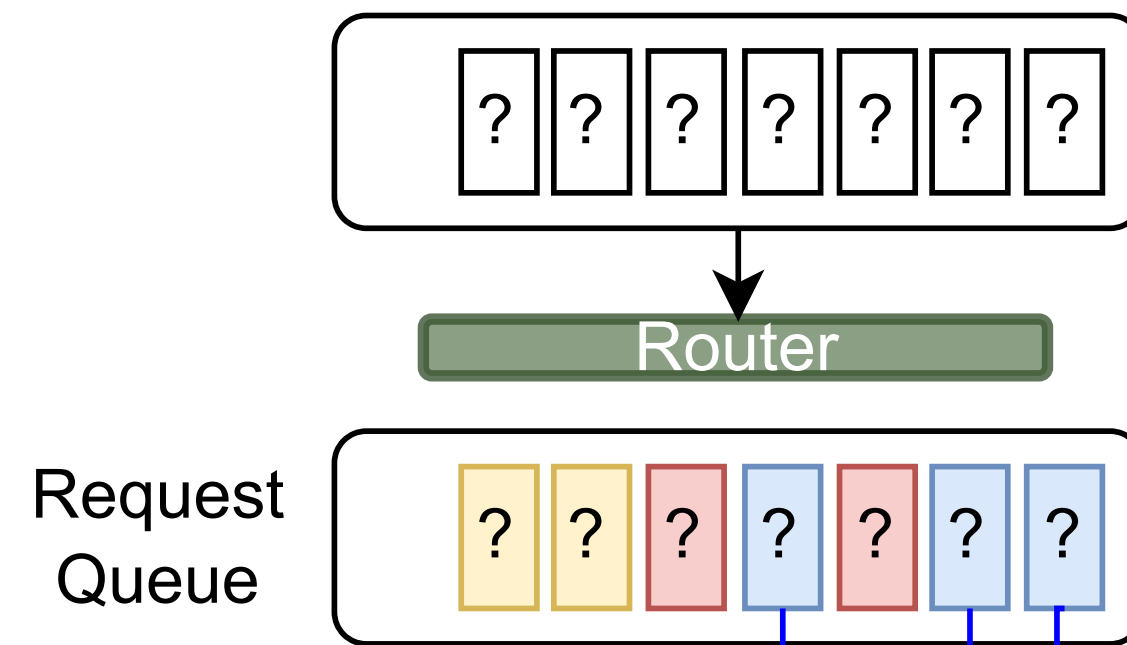
Request Queue

Pre-trained Dense Model · Refactoring · Deployment · Inference Serving at time=0 · Inference Serving at time=1

# ReadMe: Router Decoupled MoE
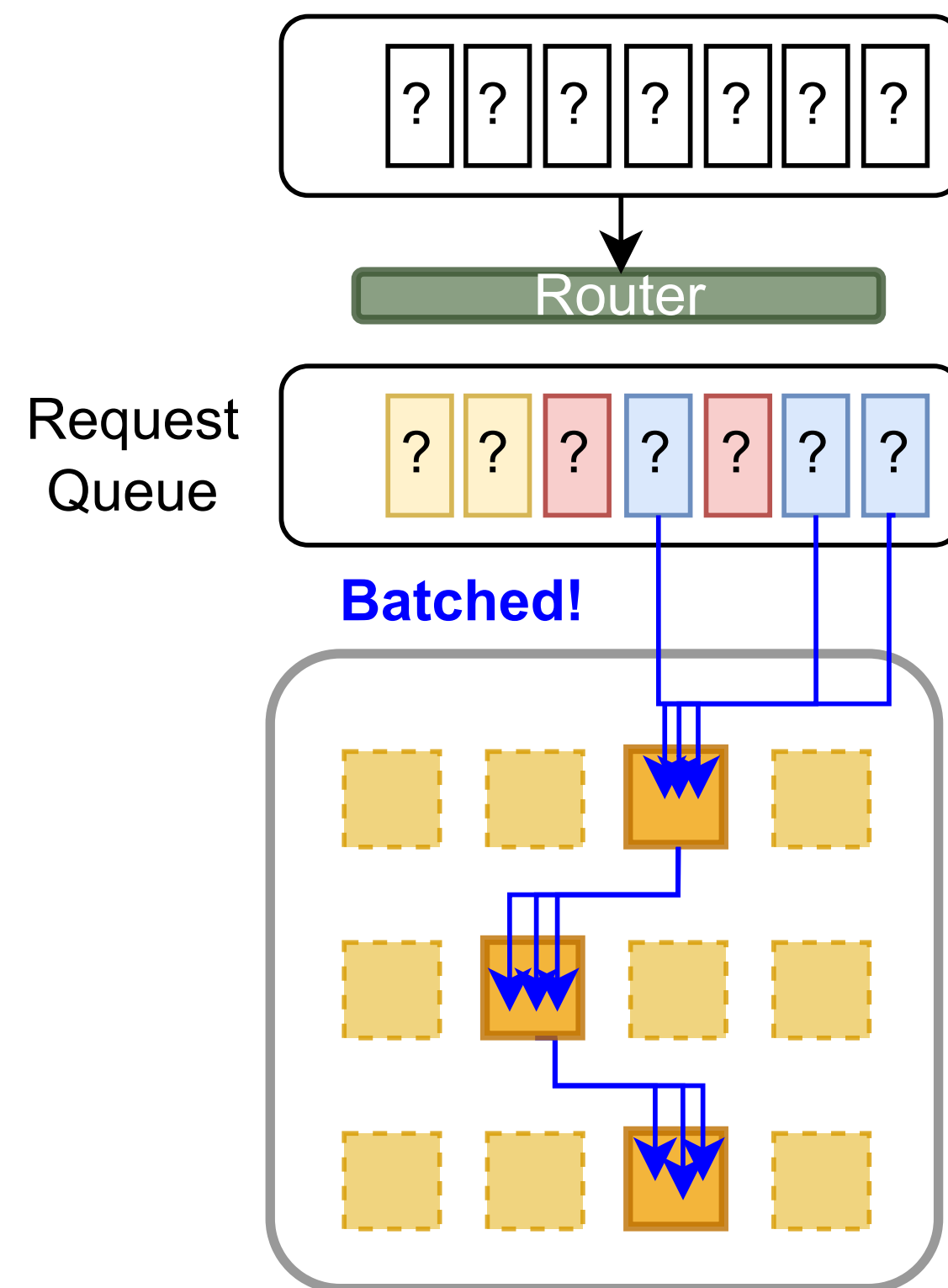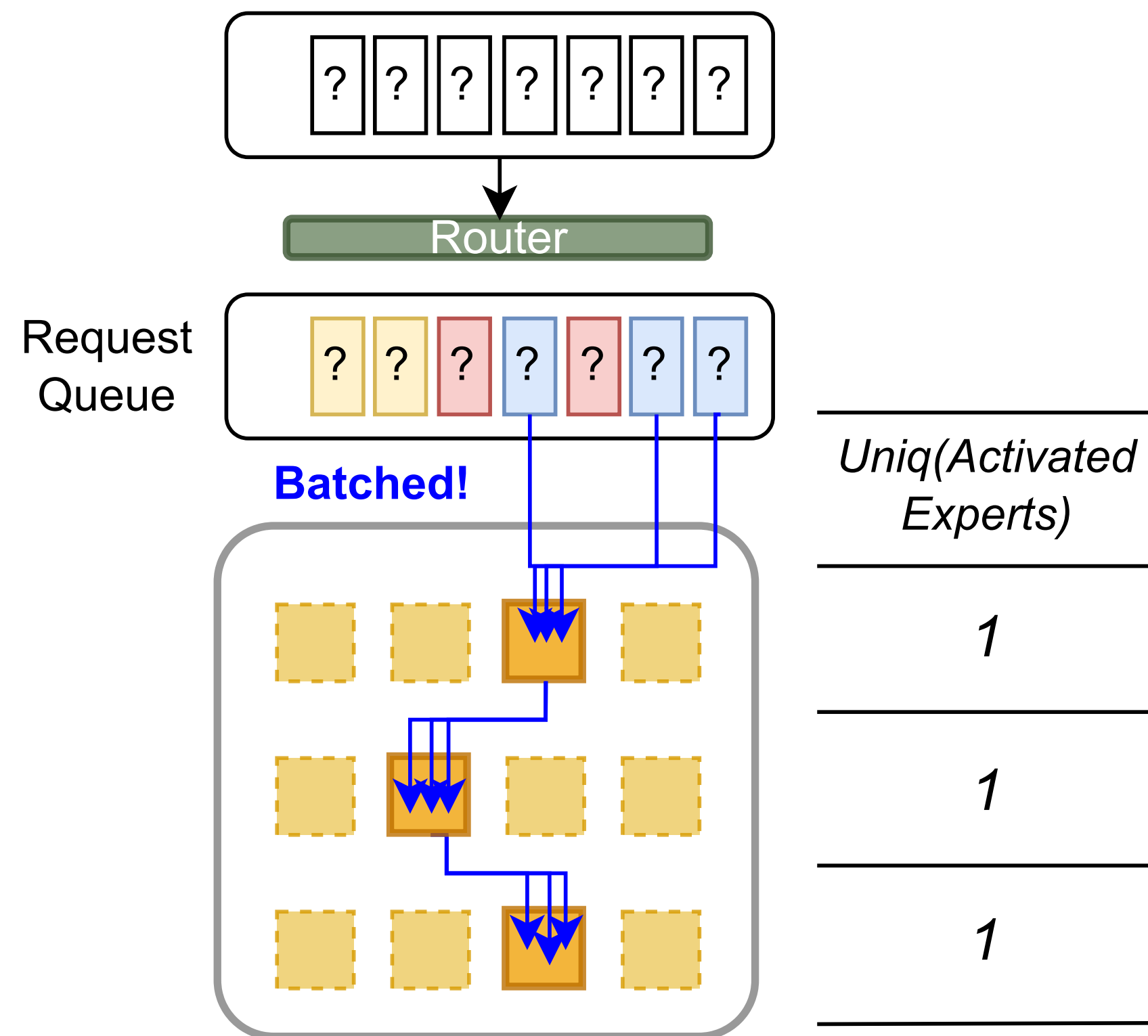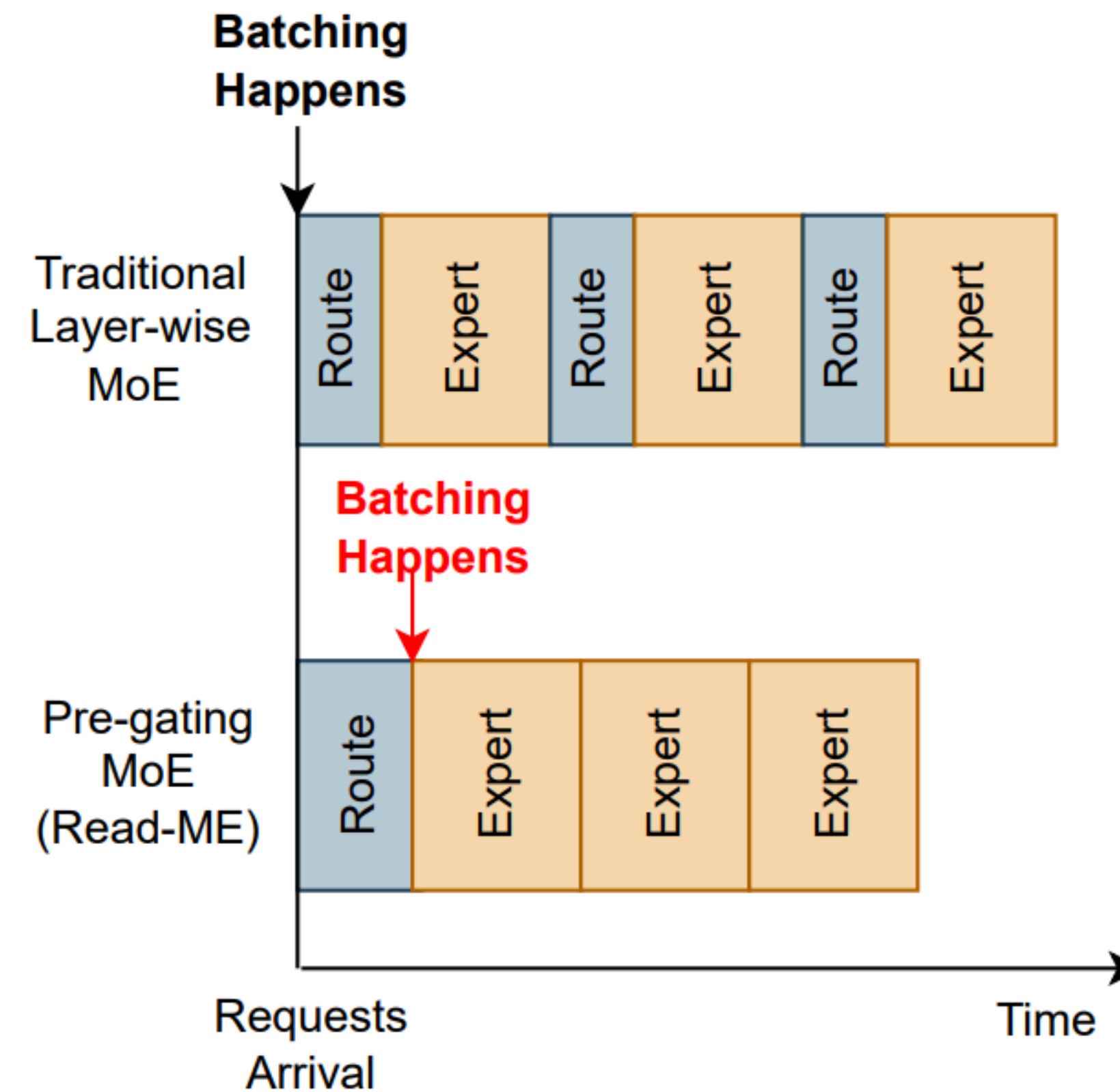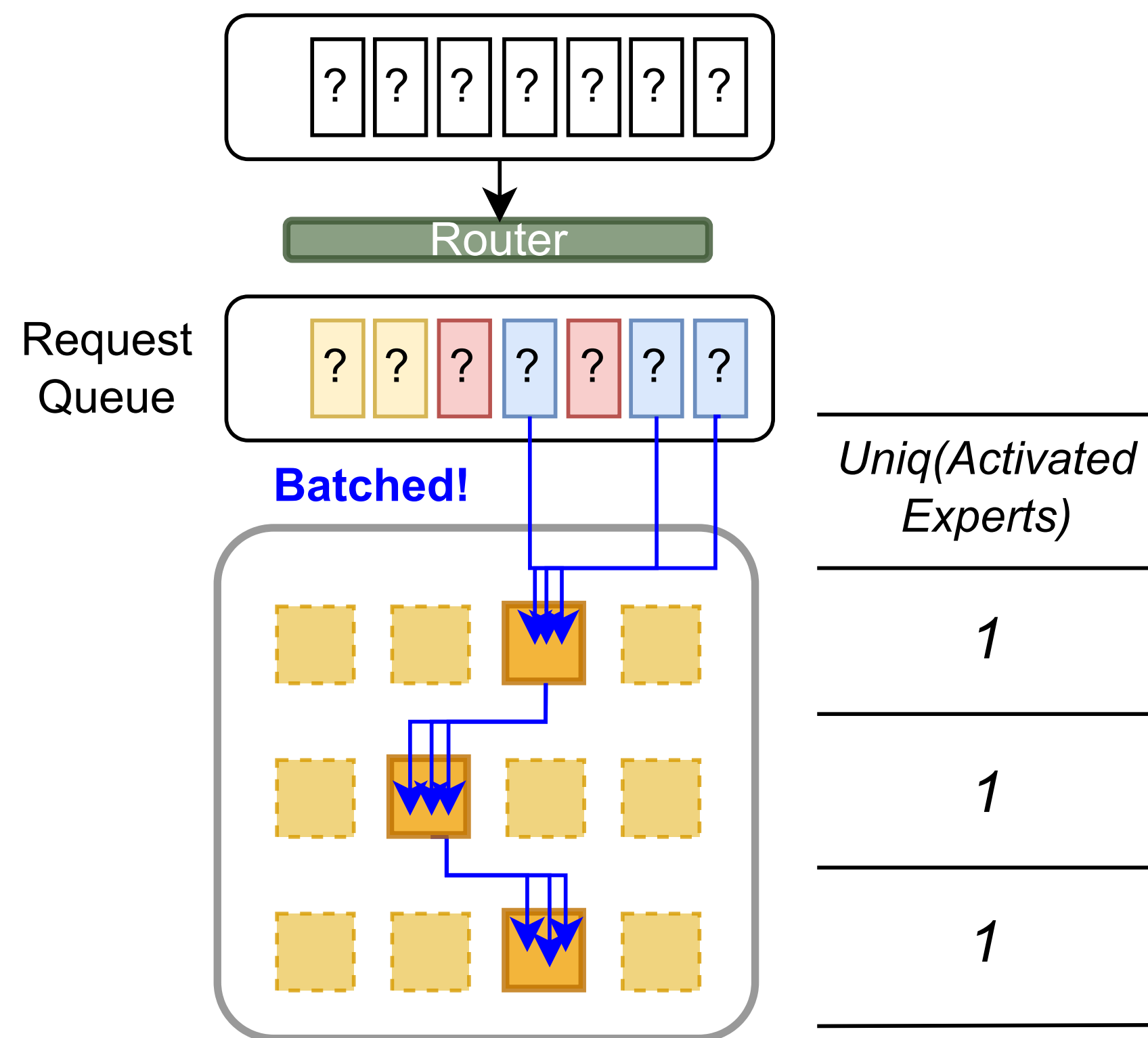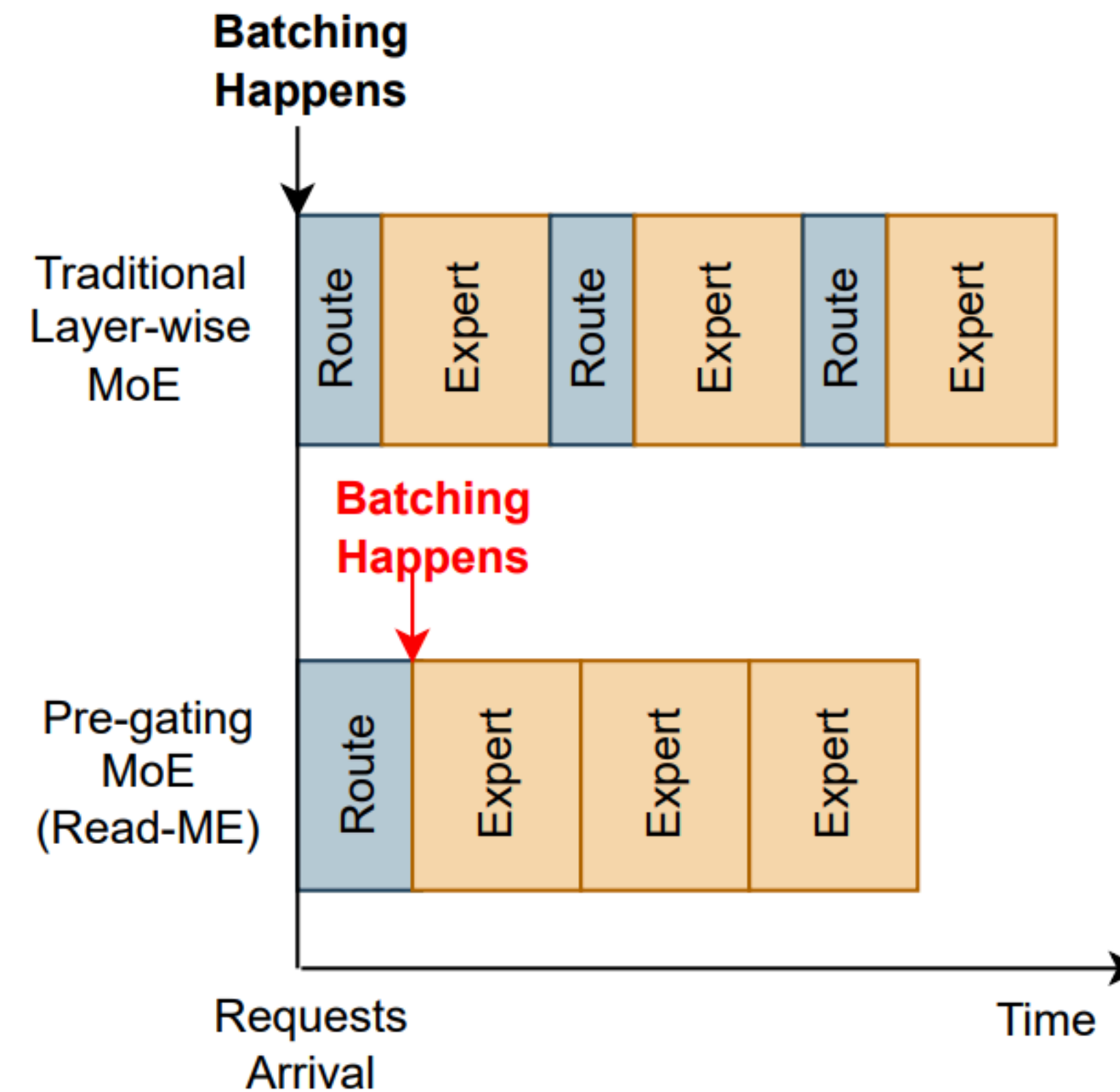## Enables Router Pre-computation and Expert-aware Batching

# ReadMe: Router Decoupled MoE
## Enables Router Pre-computation and Expert-aware Batching
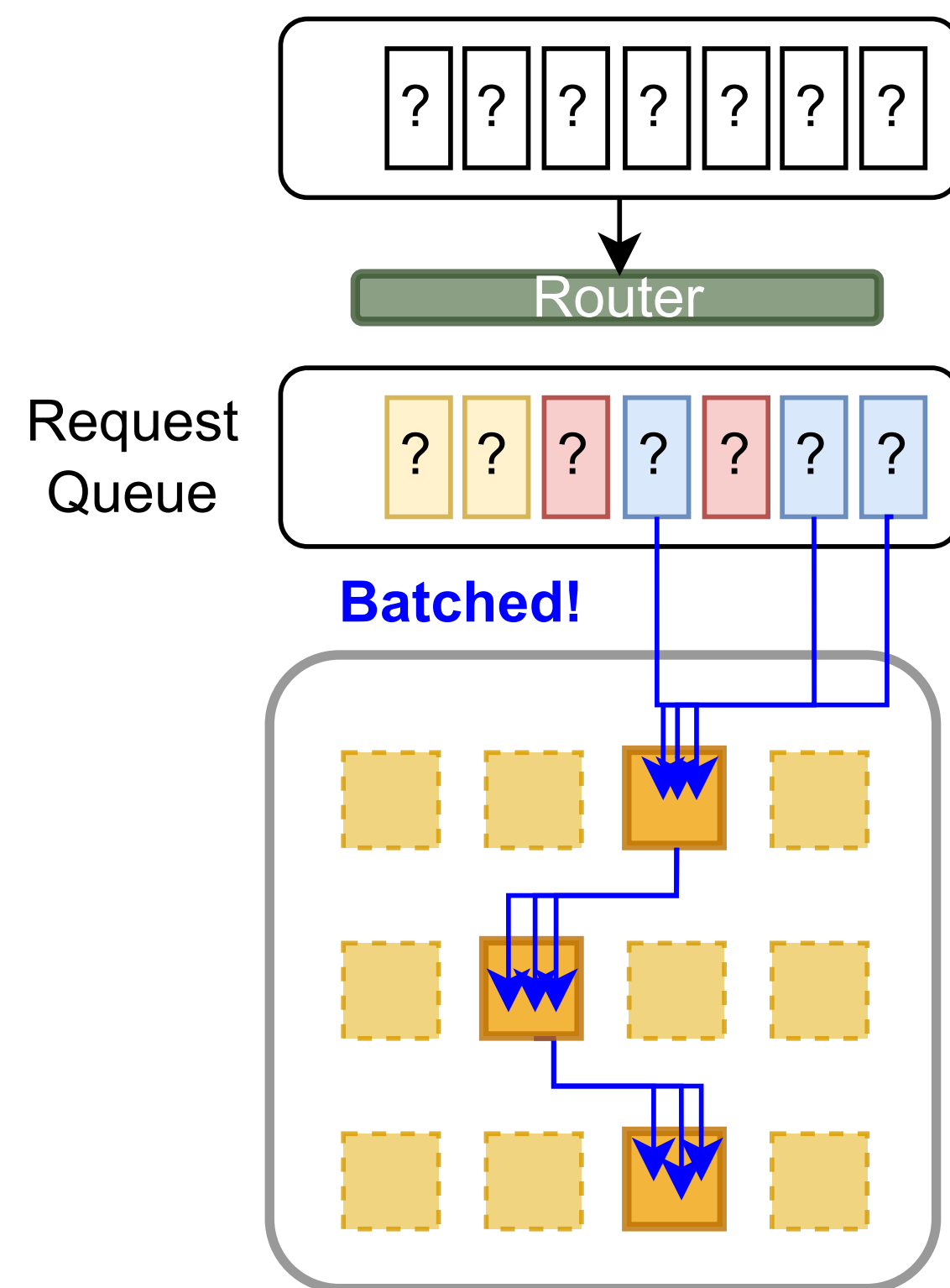
# ReadMe: Router Decoupled MoE
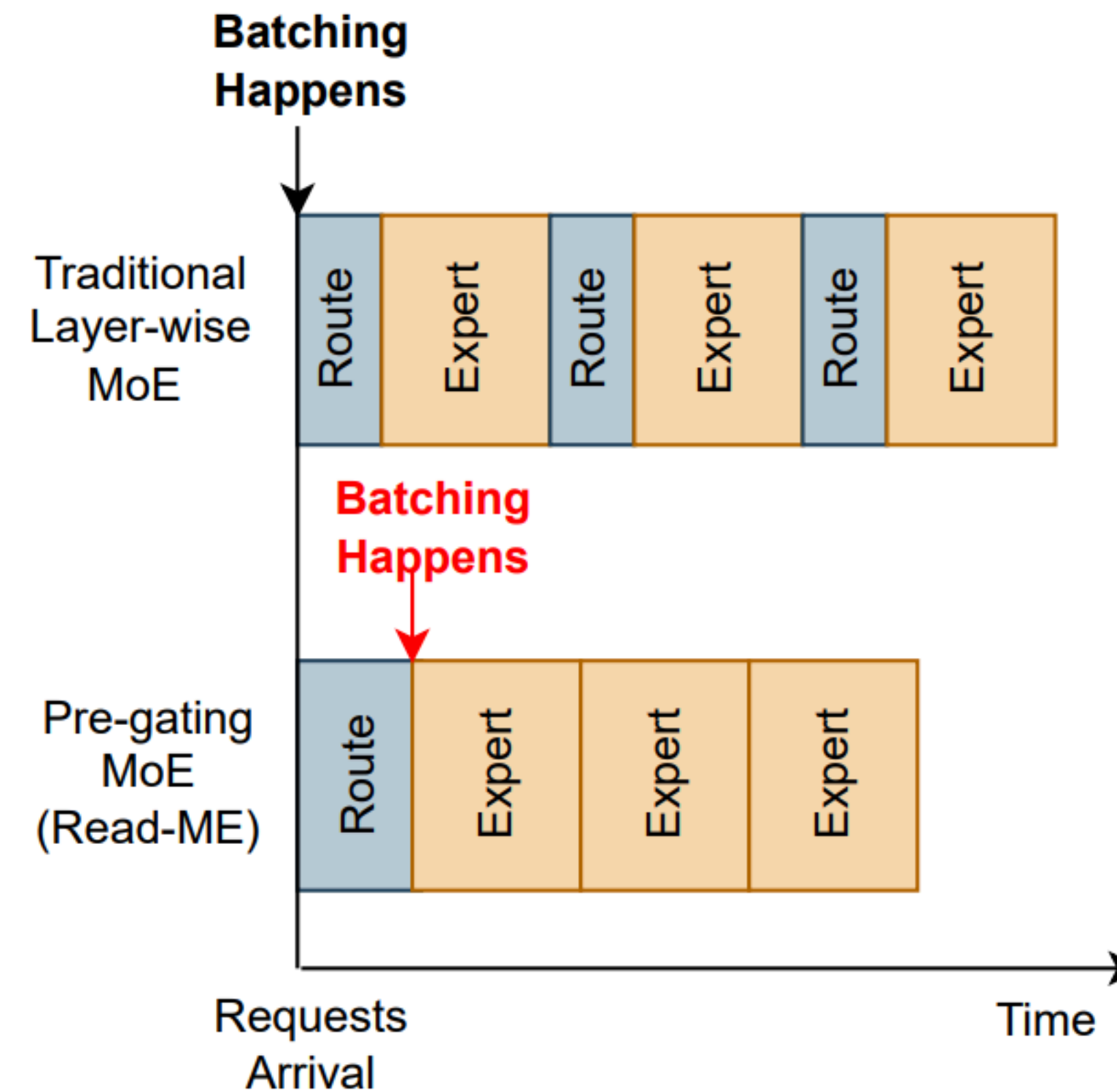## Enables Router Pre-computation and Expert-aware Batching

# ReadMe: Router Decoupled MoE
## Enables Router Pre-computation and Expert-aware Batching

# ReadMe: Router Decoupled MoE
## Enables Router Pre-computation and Expert-aware Batching
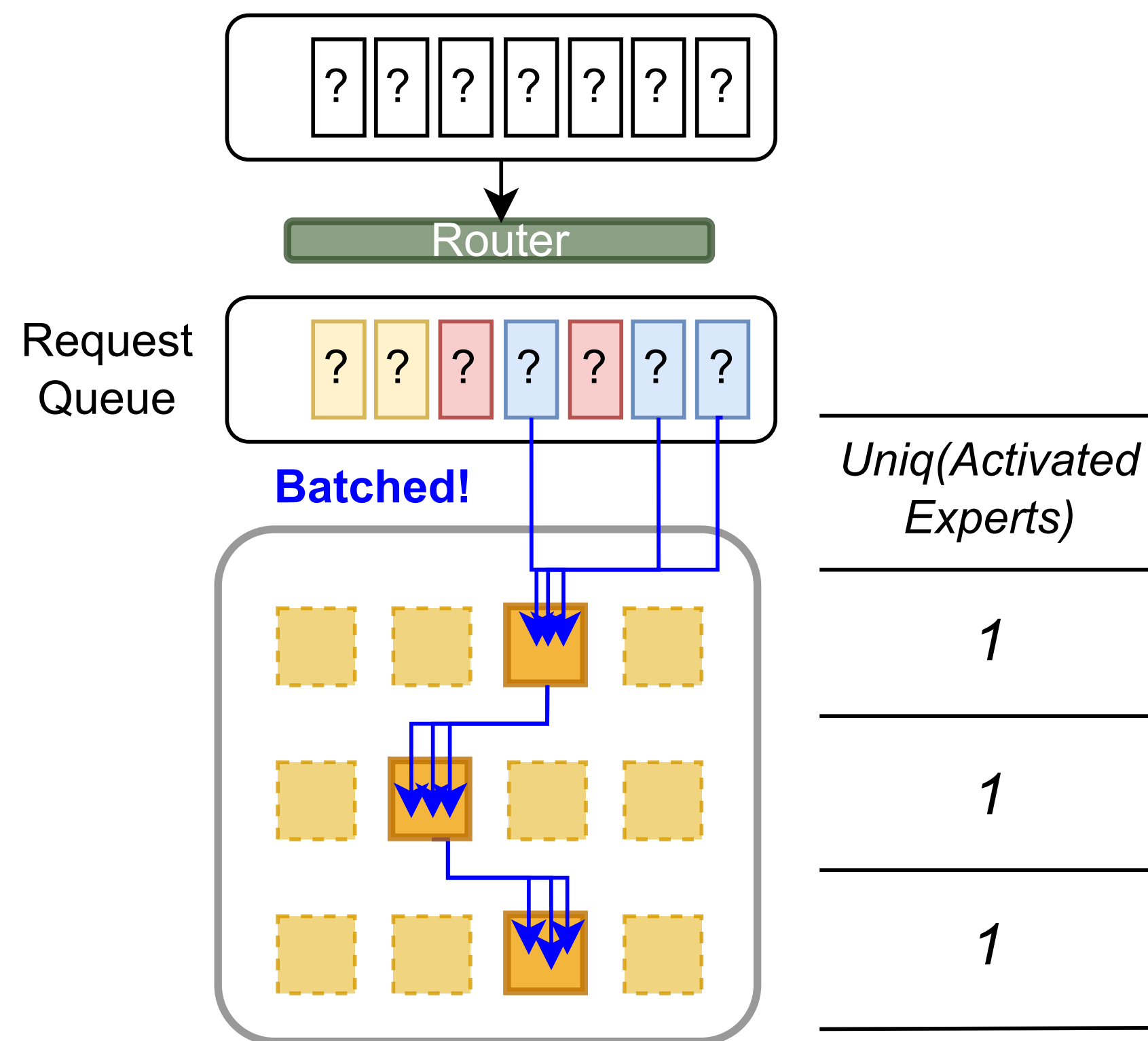
# ReadMe: Router Decoupled MoE
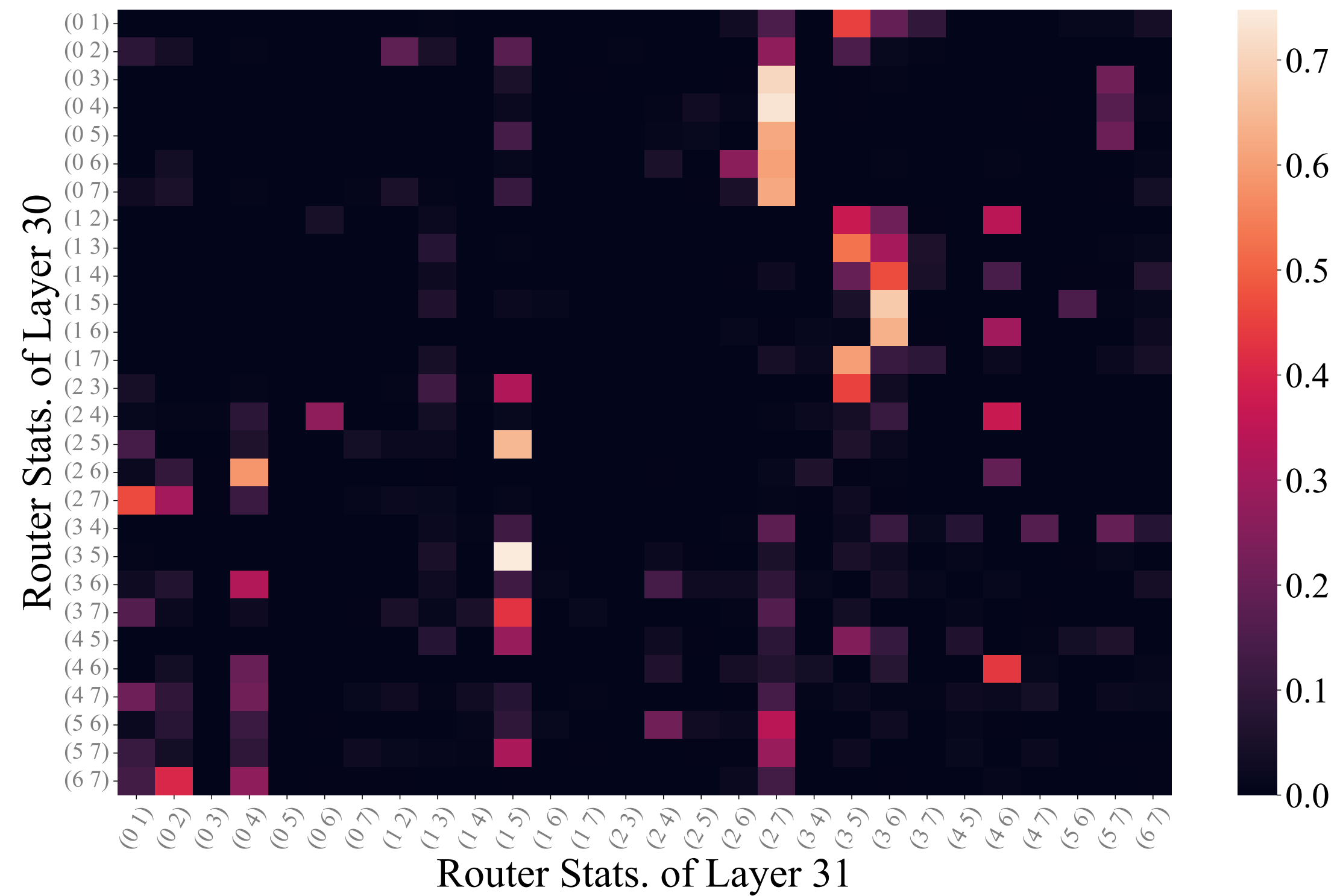## Enables Router Pre-computation and Expert-aware Batching

# ReadMe: Router Decoupled MoE
## Enables Router Pre-computation and Expert-aware Batching

# How is that possible? 😮
## The redundancy of Layer-wise Router

# Friday, 13 Dec 11AM-2PM
# Poster Session 5