

# MultiSkipGraph

A Self-stabilizing Overlay Network that Maintains Monotonic Searchability

Linghui Luo, Christian Scheideler and Thim Strothmann

Department of Computer Science  
Paderborn University

May 15, 2019

- 1 Introduction
  - Model
  - Problem Statement
- 2 Limitations
- 3 Two Protocols
  - MULTISKIPGRAPH
  - MULTISKIPGRAPH\*
- 4 Demo
- 5 Experimental Comparison
  - Comparison in Stabilization
  - Comparison in Searchability
- 6 Summary

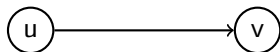
- 1 Introduction
  - Model
  - Problem Statement
- 2 Limitations
- 3 Two Protocols
  - MULTISKIPGRAPH
  - MULTISKIPGRAPH\*
- 4 Demo
- 5 Experimental Comparison
  - Comparison in Stabilization
  - Comparison in Searchability
- 6 Summary

# Assumptions

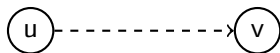
- Asynchronous distributed system
- Nodes communicate via message-passing
- Each node has:
  - a reference –  $u$
  - a unique id –  $u.id$
  - a channel –  $u.channel$
- A node  $u$  can send messages to a node  $v$  if  $u$  has the reference of  $v$
- Each node  $u$  initiates  $SEARCH(u, wID)$  messages

# Channel Connectivity Multigraph

- Graph  $G = (V, E_e \cup E_i)$ ,
  - $V$ : set of all nodes
  - $E_e$ : set of all explicit edges
  - $E_i$ : set of all implicit edges.
- Explicit Edge  $(u, v)$ : If node  $u$  stores the reference of node  $v$  locally.

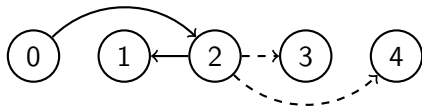


- Implicit Edge  $(u, v)$ : If there is a message in  $u.channel$  containing the reference of node  $v$ .

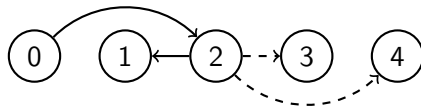


- 1 Introduction
  - Model
  - Problem Statement
- 2 Limitations
- 3 Two Protocols
  - MULTISKIPGRAPH
  - MULTISKIPGRAPH\*
- 4 Demo
- 5 Experimental Comparison
  - Comparison in Stabilization
  - Comparison in Searchability
- 6 Summary

# Topological Self-stabilization



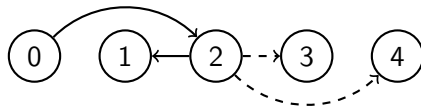
# Topological Self-stabilization



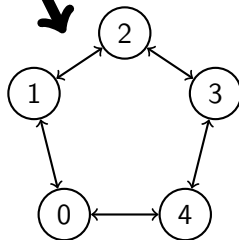
A line topology



# Topological Self-stabilization

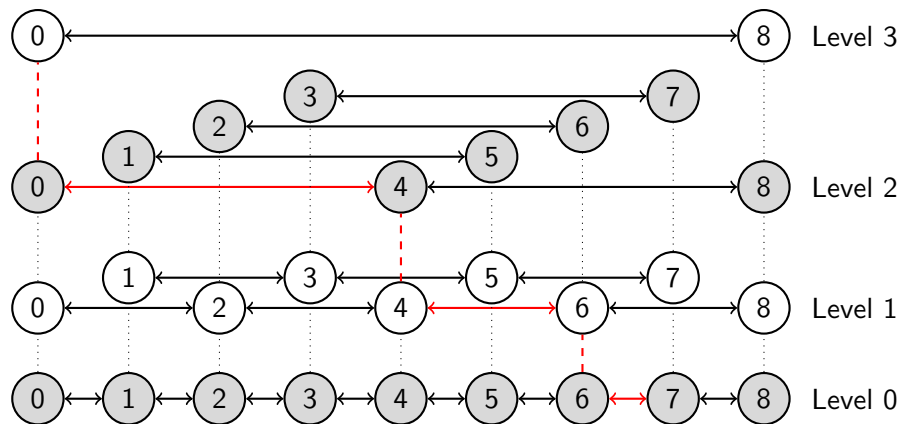


A line topology



A ring topology

# Perfect Skip Graph



Routing protocol: Greedy Search

# Monotonic Searchability

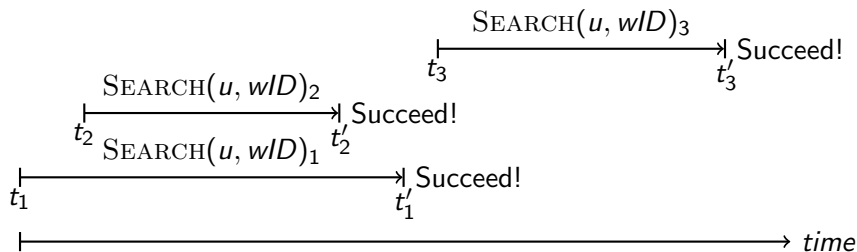


Figure: Example satisfies monotonic searchability

# Problem: Message Delivery

No protocol can **unconditionally** satisfy monotonic searchability, since messages are not necessarily delivered in FIFO order.

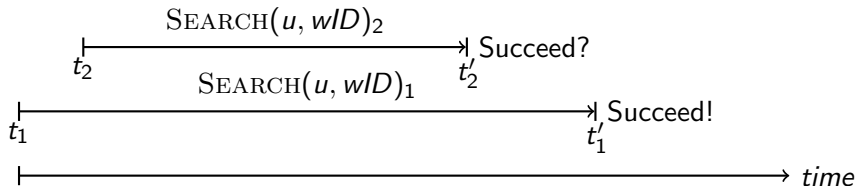


Figure: Messages initiated earlier could arrive later than messages initiated later

# Problem: Message Delivery

No protocol can **unconditionally** satisfy monotonic searchability, since messages are not necessarily delivered in FIFO order.

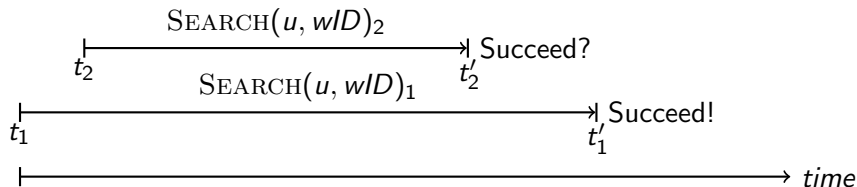


Figure: Messages initiated earlier could arrive later than messages initiated later

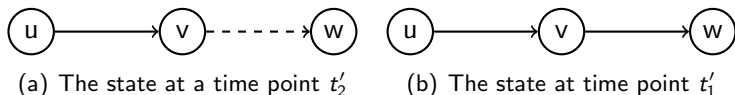
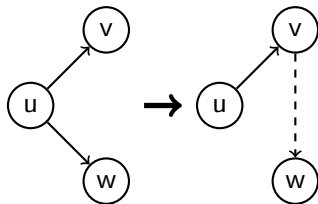


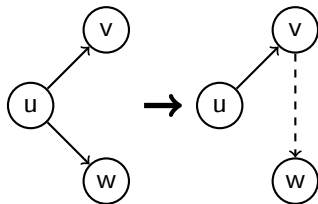
Figure: Example violates monotonic searchability

# Problem: Delegation



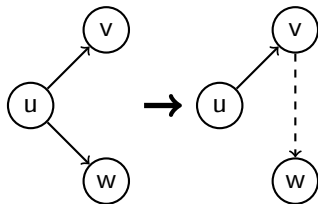
- Explicit edge  $(u, w)$  is removed in simple delegation.

# Problem: Delegation



- Explicit edge  $(u, w)$  is removed in simple delegation.
- After delegation node  $w$  is not reachable from node  $u$ .

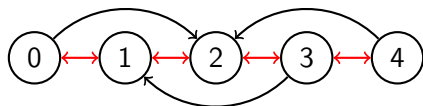
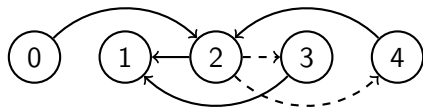
# Problem: Delegation



- Explicit edge  $(u, w)$  is removed in simple delegation.
- After delegation node  $w$  is not reachable from node  $u$ .
- A safe way: keep every explicit edge

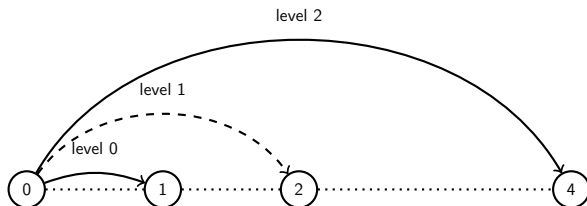


# Relaxed Self-stabilization

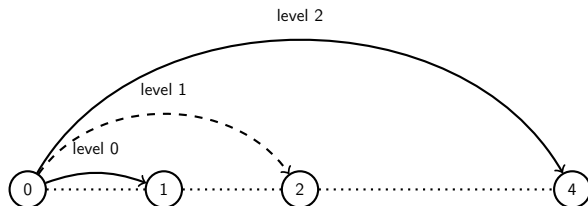


Supergraph of the line topology

# Limitations in Topology Maintenance

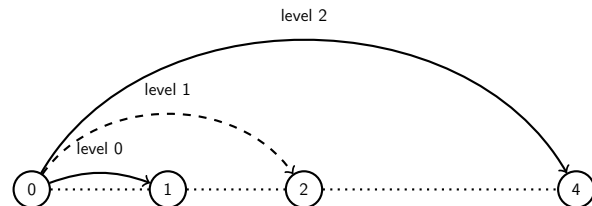


# Limitations in Topology Maintenance

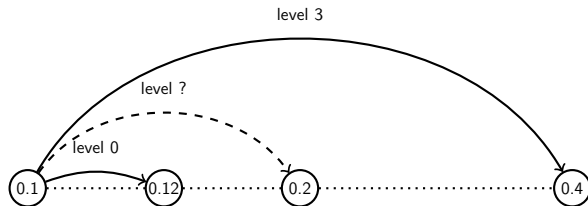


- Generally: the ids are not sequential!

# Limitations in Topology Maintenance

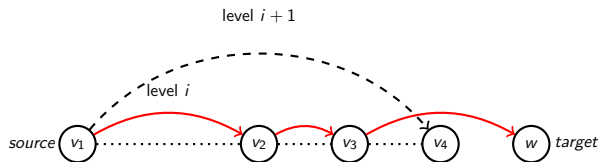


- Generally: the ids are not sequential!
- Perfect skip graph is not locally checkable.

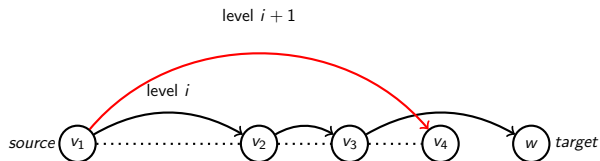


# Limitations in Monotonic Searchability

Greedy Search can not guarantee monotonic searchability



(a) First search path  $v_1 \rightarrow v_2 \rightarrow v_3 \rightarrow w$



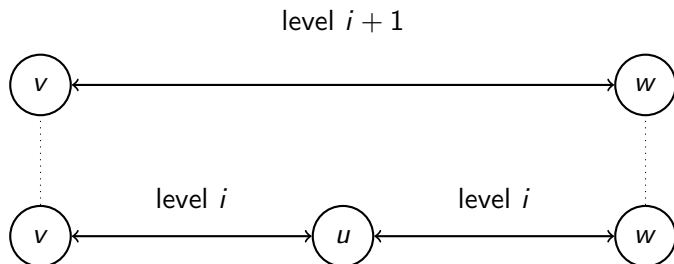
(b) Second search path  $v_1 \rightarrow v_4 \rightarrow \text{abort}$

Figure: Example violates monotonic searchability

- 1 Introduction
  - Model
  - Problem Statement
- 2 Limitations
- 3 Two Protocols
  - **MULTISKIPGRAPH**
  - MULTISKIPGRAPH\*
- 4 Demo
- 5 Experimental Comparison
  - Comparison in Stabilization
  - Comparison in Searchability
- 6 Summary

# Observation

Bottom-up process: higher levels are built on top of lower levels.



Periodically call `TIMEOUT()` action:

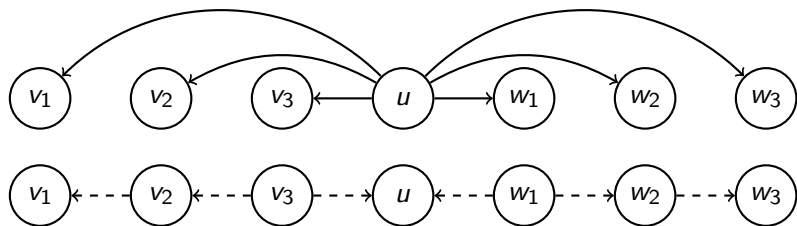
- the self-stabilizing part for topology maintenance
- the HybridSearch part for handling search requests



# MULTISKIPGRAPH

## The self-stabilizing part: linearize level-0

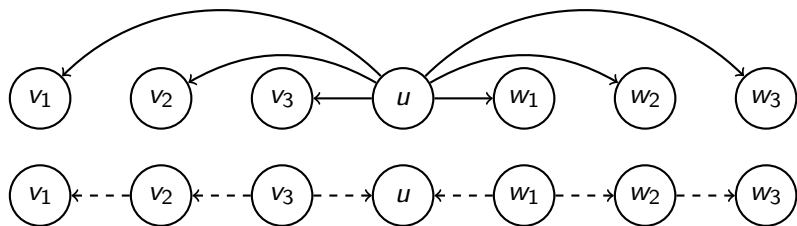
- `TIMEOUT()`: Linear introduction by sending `INTRODUCE(*)` messages



# MULTISKIPGRAPH

## The self-stabilizing part: linearize level-0

- `TIMEOUT()`: Linear introduction by sending `INTRODUCE(*)` messages

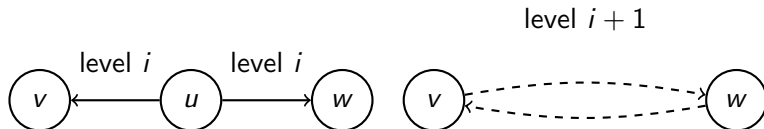


- `INTRODUCE(v)`:
  - Keep the closest neighbors as level-0 neighbors
  - Mark the old level-0 neighbors as *unknown*
  - Keep the references to *unknown* neighbors.

# MULTISKIPGRAPH

The self-stabilizing part: build level- $i + 1$  on level- $i$

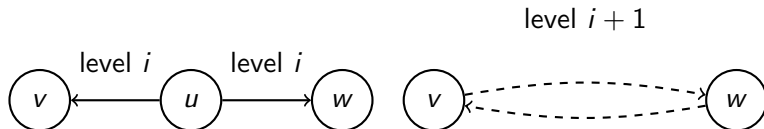
- `TIMEOUT()`: Introduce level- $i$  neighbors for level- $i + 1$  by sending `INTROLEVELNODE(*, i + 1)` messages for each level  $i$



# MULTISKIPGRAPH

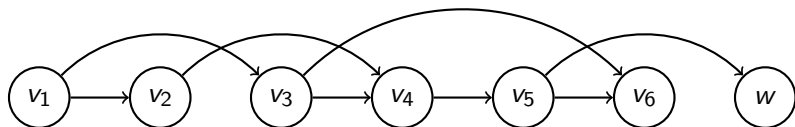
The self-stabilizing part: build level- $i + 1$  on level- $i$

- `TIMEOUT()`: Introduce level- $i$  neighbors for level- $i + 1$  by sending `INTROLEVELNODE(*, i + 1)` messages for each level  $i$

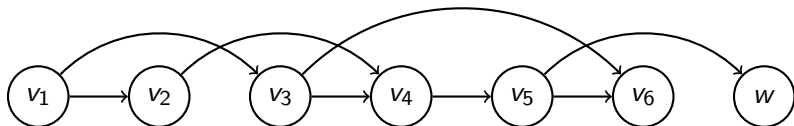


- `INTROLEVELNODE(v, i)`:
  - Always believe the newest information is correct
  - Mark old level- $i$  neighbor as *unknown*
  - Keep the references to *unknown* neighbors

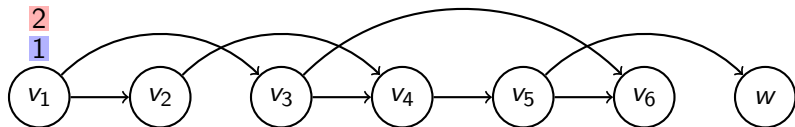
SEARCH( $v_1$ ,  $wID$ )



SEARCH( $v_1$ ,  $wID$ )

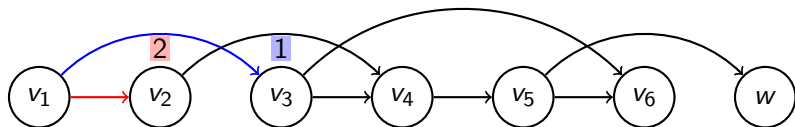


- $v_1$  doesn't forward SEARCH( $v_1$ ,  $wID$ ) immediately, but buffers it.
- TIMEOUT():  $v_1$  sends GREEDYPROBE() and GENERICPROBE() messages to itself to start two parallel probing processes.

SEARCH( $v_1$ ,  $wID$ )

- **1** GREEDYPROBE( $v_1$ ,  $wID$ ,  $seq$ )
  - Node  $v_1$  forwards it to  $v_3$
- **2** GENERICPROBE( $v_1$ ,  $wID$ ,  $Next$ ,  $seq$ ),  $Next = \{v_1\}$ 
  - Node  $v_1$  updates  $Next' = \{v_2, v_3\}$
  - Node  $v_1$  forwards GENERICPROBE( $v_1$ ,  $wID$ ,  $Next'$ ,  $seq$ ) to  $v_2$ .

SEARCH( $v_1$ ,  $wID$ )

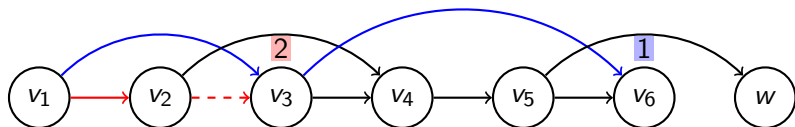


- **1** GREEDYPROBE( $v_1$ ,  $wID$ ,  $seq$ )
  - Node  $v_3$  forwards it to  $v_6$
- **2** GENERICPROBE( $v_1$ ,  $wID$ ,  $Next$ ,  $seq$ ),  $Next = \{v_2, v_3\}$ 
  - Node  $v_2$  updates  $Next' = \{v_3, v_4\}$
  - Node  $v_2$  forwards GENERICPROBE( $v_1$ ,  $wID$ ,  $Next'$ ,  $seq$ ) to  $v_3$ .



SEARCH( $v_1$ ,  $wID$ )

SEARCH( $v_1$ ,  $wID$ )



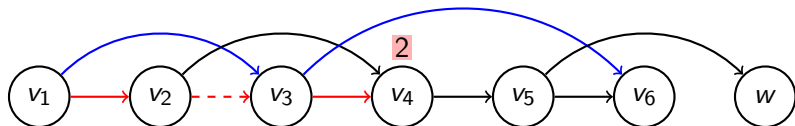
- New SEARCH( $v_1$ ,  $wID$ ) is initiated, buffering it
- **1** GREEDYPROBE( $v_1$ ,  $wID$ ,  $seq$ )
  - Node  $v_6$  can't forward it anymore, greedy probe process failed.
- **2** GENERICPROBE( $v_1$ ,  $wID$ ,  $Next$ ,  $seq$ ),  $Next = \{v_3, v_4\}$ 
  - Node  $v_3$  updates  $Next' = \{v_4, v_6\}$
  - Node  $v_3$  forwards GENERICPROBE( $v_1$ ,  $wID$ ,  $Next'$ ,  $seq$ ) to  $v_4$ .

# MULTISKIPGRAPH

## The HybridSearch part

SEARCH( $v_1$ ,  $wID$ )

SEARCH( $v_1$ ,  $wID$ )



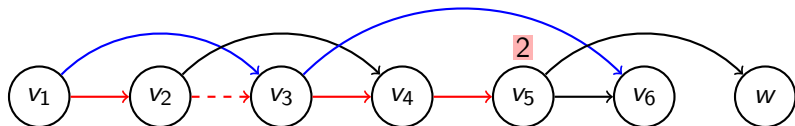
- **1** GREEDYPROBE( $v_1$ ,  $wID$ ,  $seq$ ) failed
- **2** GENERICPROBE( $v_1$ ,  $wID$ ,  $Next$ ,  $seq$ ),  $Next = \{v_4, v_6\}$ 
  - Node  $v_4$  updates  $Next' = \{v_5, v_6\}$
  - Node  $v_4$  forwards GENERICPROBE( $v_1$ ,  $wID$ ,  $Next'$ ,  $seq$ ) to  $v_5$ .

# MULTISKIPGRAPH

## The HybridSearch part

SEARCH( $v_1$ ,  $w/D$ )

SEARCH( $v_1$ ,  $w/D$ )



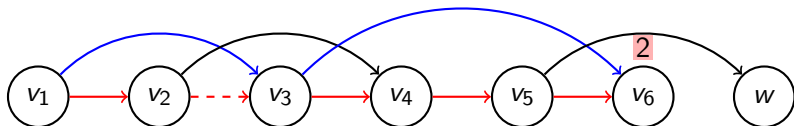
- **1** GREEDYPROBE( $v_1$ ,  $w/D$ ,  $seq$ ) failed
- **2** GENERICPROBE( $v_1$ ,  $w/D$ ,  $Next$ ,  $seq$ ),  $Next = \{v_5, v_6\}$ 
  - Node  $v_5$  updates  $Next' = \{v_6, w\}$
  - Node  $v_5$  forwards GENERICPROBE( $v_1$ ,  $w/D$ ,  $Next'$ ,  $seq$ ) to  $v_6$ .

# MULTISKIPGRAPH

## The HybridSearch part

SEARCH( $v_1$ ,  $w/D$ )

SEARCH( $v_1$ ,  $w/D$ )



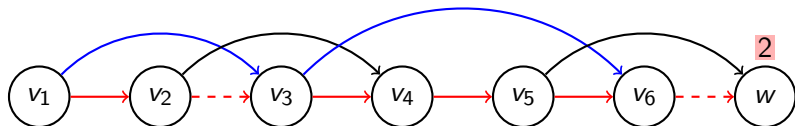
- 1 GREEDYPROBE( $v_1$ ,  $w/D$ ,  $seq$ ) failed
- 2 GENERICPROBE( $v_1$ ,  $w/D$ ,  $Next$ ,  $seq$ ),  $Next = \{v_6, w\}$ 
  - Node  $v_6$  updates  $Next' = \{w\}$
  - Node  $v_6$  forwards GENERICPROBE( $v_1$ ,  $w/D$ ,  $Next'$ ,  $seq$ ) to  $w$ .

# MULTISKIPGRAPH

## The HybridSearch part

SEARCH( $v_1$ ,  $w/D$ )

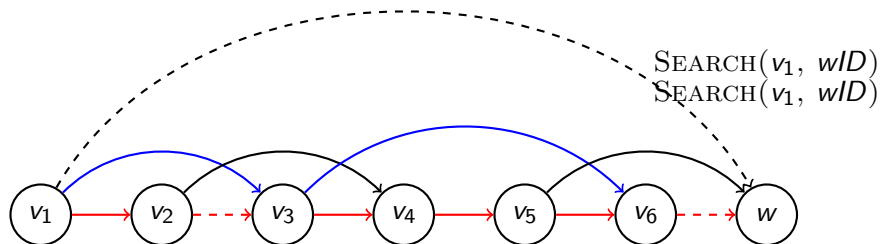
SEARCH( $v_1$ ,  $w/D$ )



- **1** GREEDYPROBE( $v_1$ ,  $w/D$ ,  $seq$ ) failed
- **2** GENERICPROBE( $v_1$ ,  $w/D$ ,  $Next$ ,  $seq$ ),  $Next = \{w\}$ 
  - Succeed!
  - Node  $w$  sends PROBESUCCESS( $w/D$ ,  $seq$ ,  $w$ ) back to  $v_1$ .

# MULTISKIPGRAPH

## The HybridSearch part



- 1 GREEDYPROBE( $v_1, wID, seq$ ) failed
- 2 GENERICPROBE( $v_1, wID, Next, seq$ ),  $Next = \{w\}$ 
  - Succeed! Node  $w$  sends PROBESUCCESS( $wID, seq, w$ ) back to  $v_1$ .
  - Node  $v_1$  sends all  $\text{SEARCH}(v_1, wID)$  to  $w$ .

How is monotonic searchability maintained?

How is monotonic searchability maintained?

- $\text{SEARCH}(v, w/D)$  messages are only sent to target node  $w$  when  $v$  receives  $\text{PROBESUCCESS}(w/D, seq, w)$  message.



How is monotonic searchability maintained?

- $\text{SEARCH}(v, w/D)$  messages are only sent to target node  $w$  when  $v$  receives  $\text{PROBESUCCESS}(w/D, seq, w)$  message.
- Node  $v$  receives  $\text{PROBESUCCESS}(w/D, seq, w)$  messages: there is a directed path from  $v$  to  $w$ .

How is monotonic searchability maintained?

- $\text{SEARCH}(v, w/D)$  messages are only sent to target node  $w$  when  $v$  receives  $\text{PROBESUCCESS}(w/D, seq, w)$  message.
- Node  $v$  receives  $\text{PROBESUCCESS}(w/D, seq, w)$  messages: there is a directed path from  $v$  to  $w$ .
- **MULTISKIPGRAPH** removes no explicit edge: the directed path from  $v$  to  $w$  always exists.

How is monotonic searchability maintained?

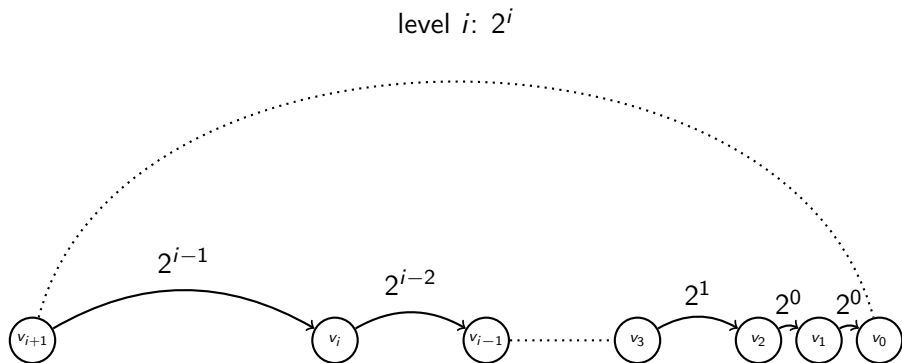
- $\text{SEARCH}(v, w/D)$  messages are only sent to target node  $w$  when  $v$  receives  $\text{PROBESUCCESS}(w/D, seq, w)$  message.
- Node  $v$  receives  $\text{PROBESUCCESS}(w/D, seq, w)$  messages: there is a directed path from  $v$  to  $w$ .
- **MULTISKIPGRAPH** removes no explicit edge: the directed path from  $v$  to  $w$  always exists.
- Probing for  $\text{SEARCH}(v, w/D)$  will always succeed.

- 1 Introduction
  - Model
  - Problem Statement
- 2 Limitations
- 3 Two Protocols
  - MULTISKIPGRAPH
  - **MULTISKIPGRAPH\***
- 4 Demo
- 5 Experimental Comparison
  - Comparison in Stabilization
  - Comparison in Searchability
- 6 Summary

# Observation

## Deterministic Search Path

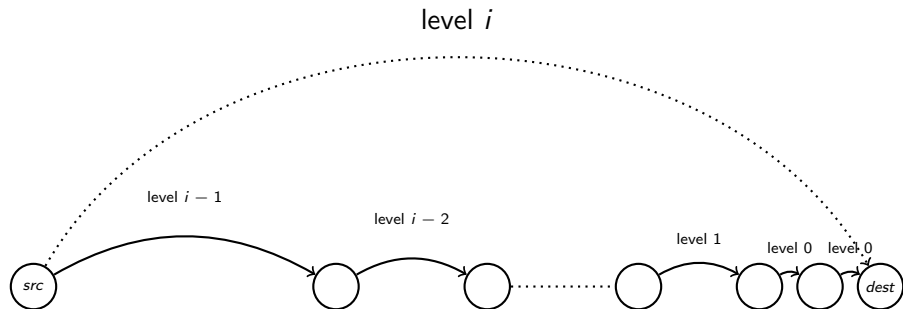
Nodes  $v_{i+1}$  and  $v_0$  are level- $i$  neighbors



# MULTISKIPGRAPH\*

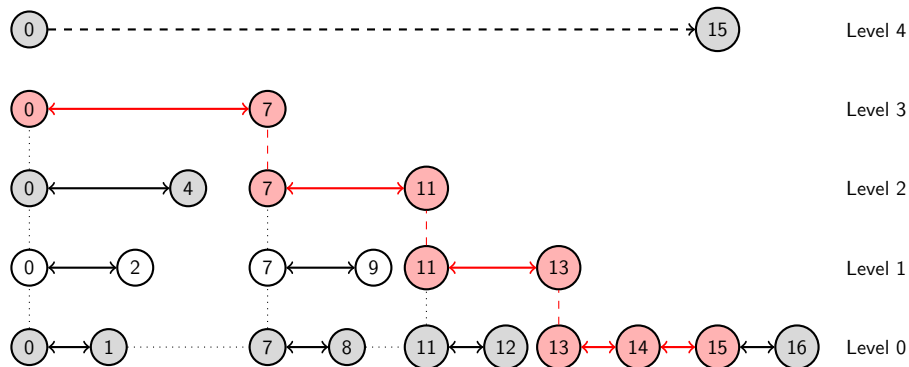
Changes in the self-stabilizing part

- More strict in updating level neighbors
- $\text{INTROLEVELNODE}(dest, i)$ : Probing the deterministic search path by forwarding  $\text{PROBELEVELNODE}(src, dest, i, j)$  message



# Observation

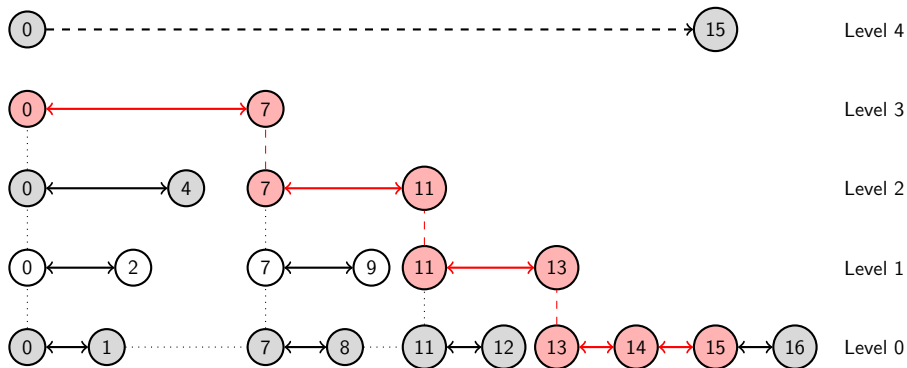
## Deterministic Search Path



- Deterministic Search Path exists, but nodes 1 and 15 are not level-4 neighbors

# Observation

## Deterministic Search Path



- Deterministic Search Path exists, but nodes 1 and 15 are not level-4 neighbors
- Still needs to update every time when probe succeeds



# MULTISKIPGRAPH\*

Changes in the self-stabilizing part

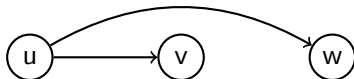
TIMEOUT(): Safely delegate the *unknown* neighbors

# MULTISKIPGRAPH\*

Changes in the self-stabilizing part

TIMEOUT(): Safely delegate the *unknown* neighbors

- $u$  wants to delegate its right *unknown* neighbor  $w$ .

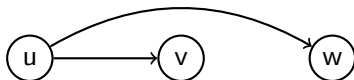


# MULTISKIPGRAPH\*

Changes in the self-stabilizing part

TIMEOUT(): Safely delegate the *unknown* neighbors

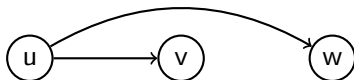
- $u$  wants to delegate its right *unknown* neighbor  $w$ .



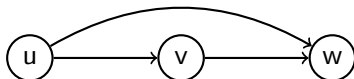
- $u$  sends  $\text{SAFEINTRODUCE}(w, u)$  to its right neighbor  $v$  which is closest to  $w$  and  $v.id < w.id$ .

TIMEOUT(): Safely delegate the *unknown* neighbors

- $u$  wants to delegate its right *unknown* neighbor  $w$ .

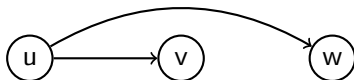


- $u$  sends  $\text{SAFEINTRODUCE}(w, u)$  to its right neighbor  $v$  which is closest to  $w$  and  $v.id < w.id$ .
- $v$  adds  $w$  as its *unknown* neighbor and sends  $\text{SAFEDELETION}(w)$  back to  $u$ .

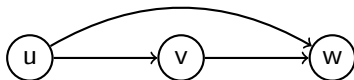


TIMEOUT(): Safely delegate the *unknown* neighbors

- $u$  wants to delegate its right *unknown* neighbor  $w$ .

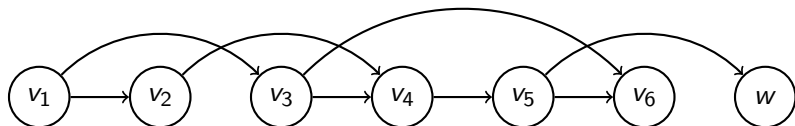


- $u$  sends  $\text{SAFEINTRODUCE}(w, u)$  to its right neighbor  $v$  which is closest to  $w$  and  $v.id < w.id$ .
- $v$  adds  $w$  as its *unknown* neighbor and sends  $\text{SAFEDELETION}(w)$  back to  $u$ .

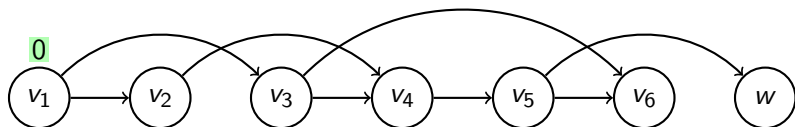


- $u$  deletes  $w$ .

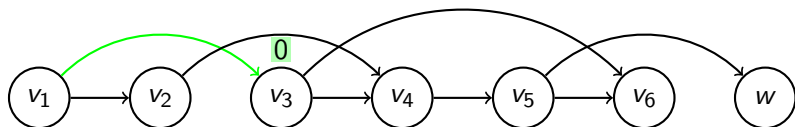


SEARCH( $v_1$ ,  $w/D$ )

- **TIMEOUT()**:  $v_1$  sends a **SLOWGREEDYPROBE()** message to itself to start a probing process.

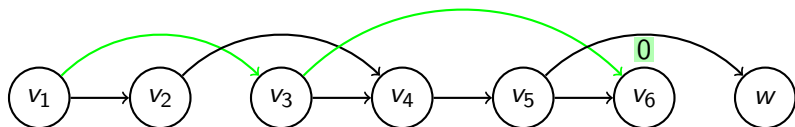
SEARCH( $v_1$ ,  $wID$ )

- **0** SLOWGREEDYPROBE( $v_1$ ,  $wID$ ,  $Prev$ ,  $Next$ ,  $seq$ ),  
 $Prev = \emptyset$ ,  $Next = \{v_1\}$ 
  - Node  $v_1$  updates  $Prev' = \{v_1\}$  and  $Next' = \{v_2, v_3\}$
  - Node  $v_1$  forwards SLOWGREEDYPROBE( $v_1$ ,  $wID$ ,  $Prev'$ ,  $Next'$ ,  $seq$ ) to  $v_3$

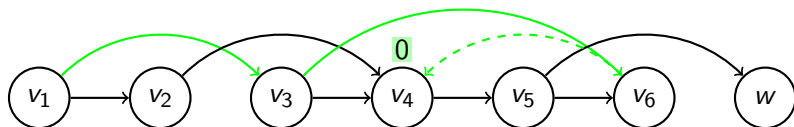
SEARCH( $v_1$ ,  $wID$ )

- **0** SLOWGREEDYPROBE( $v_1$ ,  $wID$ ,  $Prev$ ,  $Next$ ,  $seq$ ),  
 $Prev = \{v_1\}$ ,  $Next = \{v_2, v_3\}$ 
  - Node  $v_3$  updates  $Prev' = \{v_1, v_3\}$  and  $Next' = \{v_2, v_4, v_6\}$
  - Node  $v_3$  forwards SLOWGREEDYPROBE( $v_1$ ,  $wID$ ,  $Prev'$ ,  $Next'$ ,  $seq$ ) to  $v_6$

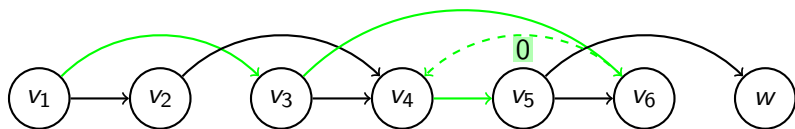


SEARCH( $v_1$ ,  $wID$ )

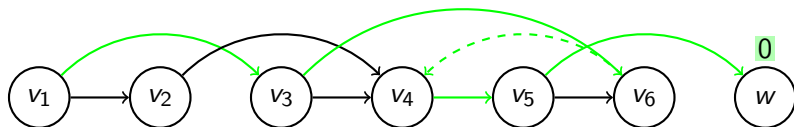
- **0** SLOWGREEDYPROBE( $v_1$ ,  $wID$ ,  $Prev$ ,  $Next$ ,  $seq$ ),  
 $Prev = \{v_1, v_3\}$ ,  $Next = \{v_2, v_4, v_6\}$ 
  - Node  $v_6$  updates  $Prev' = \{v_1, v_3, v_6\}$  and  $Next' = \{v_2, v_4\}$
  - Node  $v_6$  forwards SLOWGREEDYPROBE( $v_1$ ,  $wID$ ,  $Prev'$ ,  $Next'$ ,  $seq$ ) to  $v_4$

SEARCH( $v_1$ ,  $wID$ )

- **0** SLOWGREEDYPROBE( $v_1$ ,  $wID$ ,  $Prev$ ,  $Next$ ,  $seq$ ),  
 $Prev = \{v_1, v_3, v_6\}$ ,  $Next = \{v_2, v_4\}$ 
  - Node  $v_4$  updates  $Prev' = \{v_1, v_3, v_6, v_4\}$  and  $Next' = \{v_2, v_5\}$
  - Node  $v_4$  forwards SLOWGREEDYPROBE( $v_1$ ,  $wID$ ,  $Prev'$ ,  $Next'$ ,  $seq$ ) to  $v_5$

SEARCH( $v_1$ ,  $wID$ )

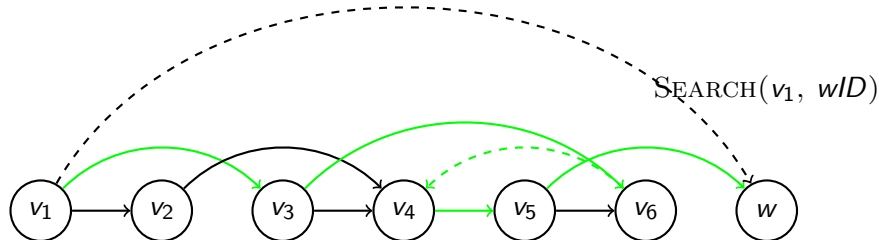
- **0** SLOWGREEDYPROBE( $v_1$ ,  $wID$ ,  $Prev$ ,  $Next$ ,  $seq$ ),  
 $Prev = \{v_1, v_3, v_6, v_4\}$ ,  $Next = \{v_2, v_5\}$ 
  - Node  $v_5$  updates  $Prev' = \{v_1, v_3, v_6, v_4, v_5\}$  and  $Next' = \{v_2, w\}$
  - Node  $v_5$  forwards SLOWGREEDYPROBE( $v_1$ ,  $wID$ ,  $Prev'$ ,  $Next'$ ,  $seq$ ) to  $w$

SEARCH( $v_1$ ,  $wID$ )

- **0** SLOWGREEDYPROBE( $v_1$ ,  $wID$ ,  $Prev$ ,  $Next$ ,  $seq$ ),  
 $Prev = \{v_1, v_3, v_6, v_4, v_5\}$ ,  $Next = \{v_2, w\}$ 
  - Succeed!
  - Node  $w$  sends PROBESUCCESS( $wID$ ,  $seq$ ,  $w$ ) back to  $v_1$ .

# MULTISKIPGRAPH\*

## The SlowGreedySearch part



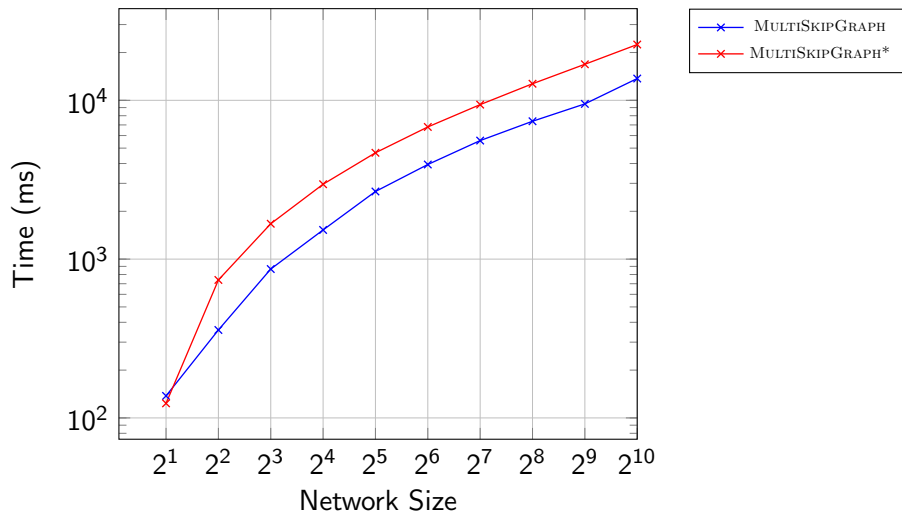
- **0**  $\text{SLOWGREEDYPROBE}(v_1, wID, \text{Prev}, \text{Next}, \text{seq})$ ,  
 $\text{Prev} = \{v_1, v_3, v_6, v_4, v_5\}$ ,  $\text{Next} = \{v_2, w\}$ 
  - Succeed! Node  $w$  sends  $\text{PROBESUCCESS}(wID, \text{seq}, w)$  back to  $v_1$ .
  - Node  $v_1$  sends  $\text{SEARCH}(v_1, wID)$  to  $w$ .

<https://www.youtube.com/watch?v=S8yd7fApSfk&list=PLSdezCrzv5YPpamQYw46HK-uNdyhyQVbt>

- 1 Introduction
  - Model
  - Problem Statement
- 2 Limitations
- 3 Two Protocols
  - MULTISKIPGRAPH
  - MULTISKIPGRAPH\*
- 4 Demo
- 5 Experimental Comparison**
  - Comparison in Stabilization
  - Comparison in Searchability
- 6 Summary

# Comparison in Stabilization

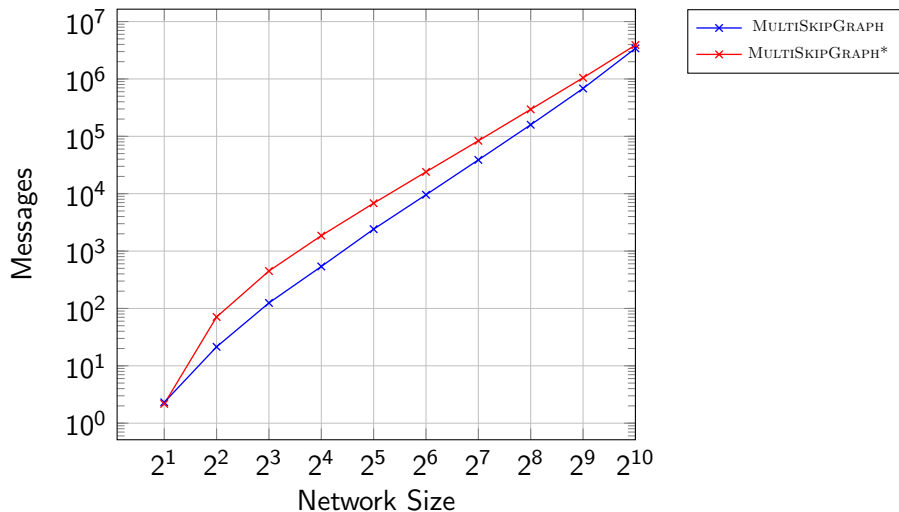
## Stabilization Time





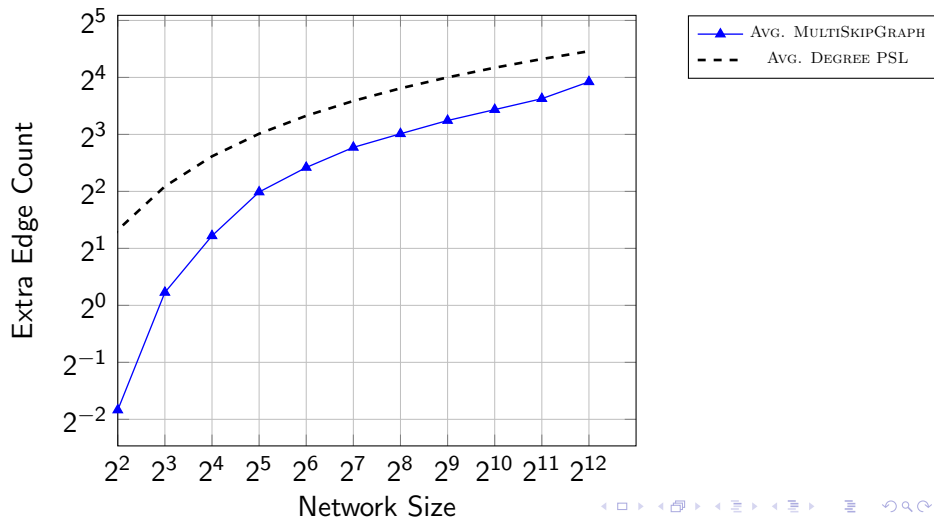
# Comparison in Stabilization

## Used Messages



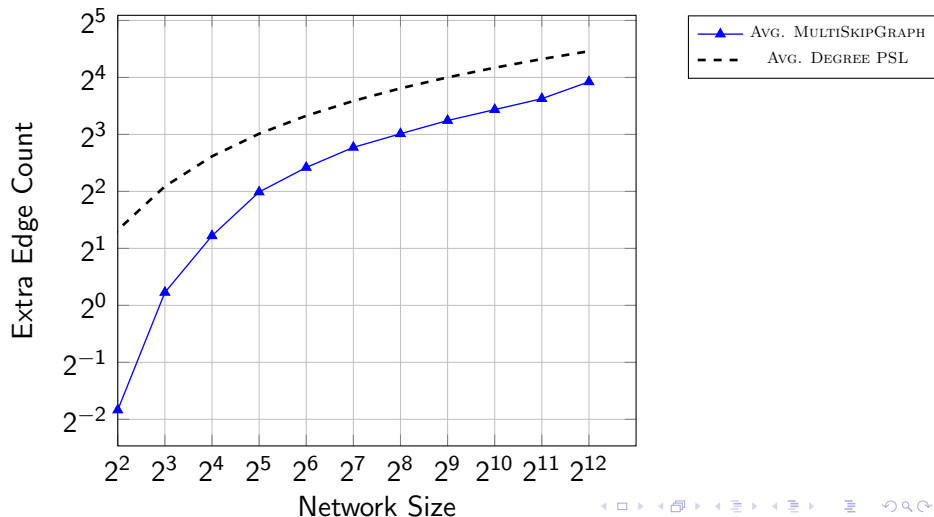
# Comparison in Stabilization

## Extra Edge Count

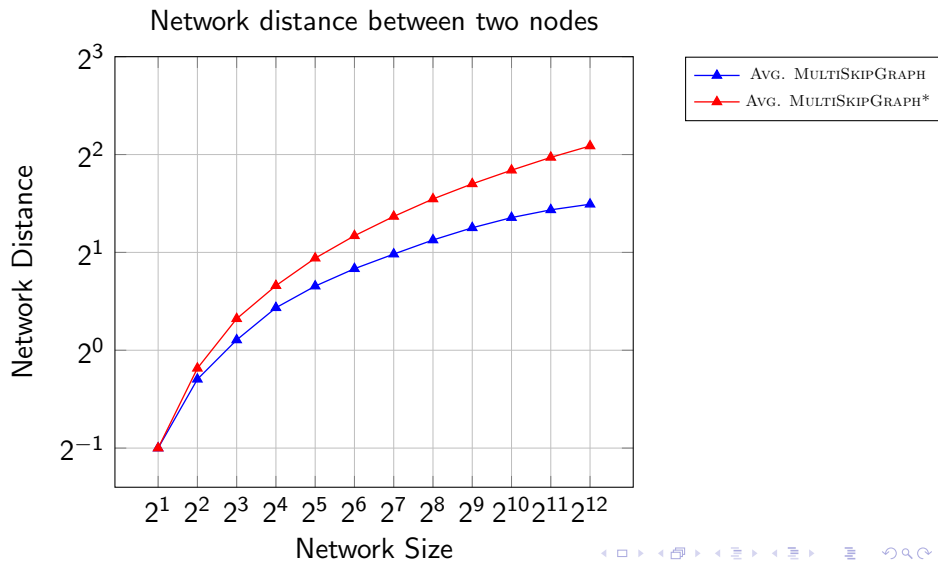


# Comparison in Stabilization

## Extra Edge Count

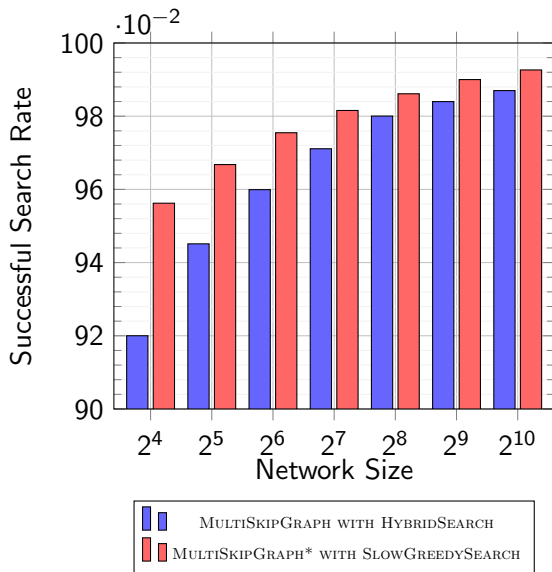


# Comparison in Stabilization

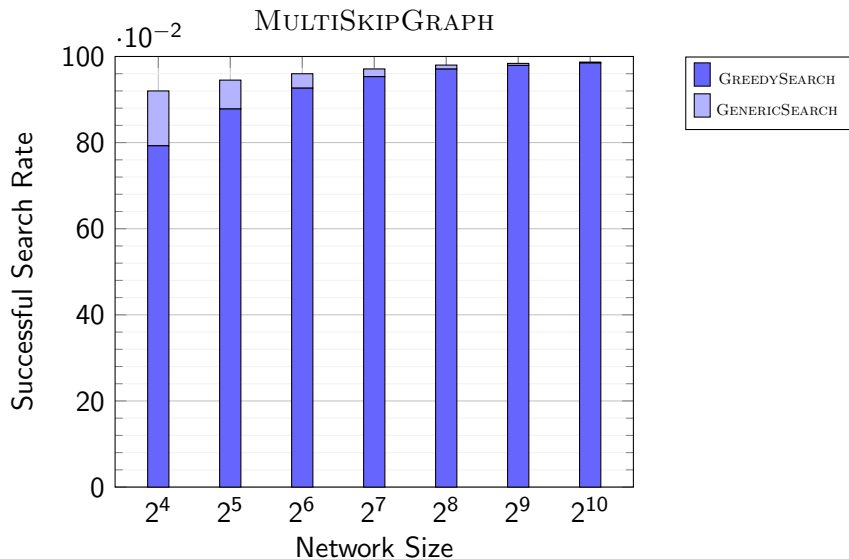


- 1 Introduction
  - Model
  - Problem Statement
- 2 Limitations
- 3 Two Protocols
  - MULTISKIPGRAPH
  - MULTISKIPGRAPH\*
- 4 Demo
- 5 Experimental Comparison**
  - Comparison in Stabilization
  - **Comparison in Searchability**
- 6 Summary

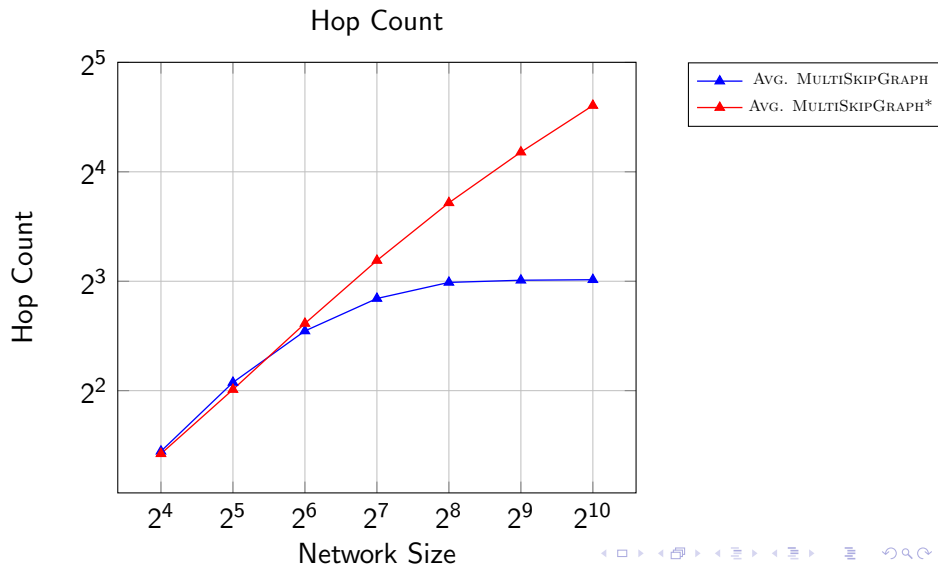
# Comparison in Searchability



# Comparison in Searchability



# Comparison in Searchability





- MULTISKIPGRAPH
  - Fast
  - Final topology: supergraph of the perfect skip graph
  - Extra local memory overhead
  - Reduced search time
- MULTISKIPGRAPH\*
  - Final topology: exactly the perfect skip graph
  - More messages
  - Higher successful search rate with SlowGreedySearch
- Outlook
  - Churn management – node insertion and deletion
  - Extend to circular skip graph for more robustness
  - Experiments on multiple computers

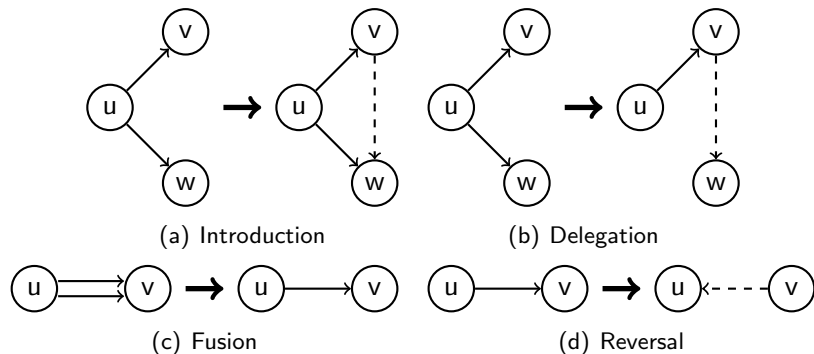


Figure: Illustration of four primitives

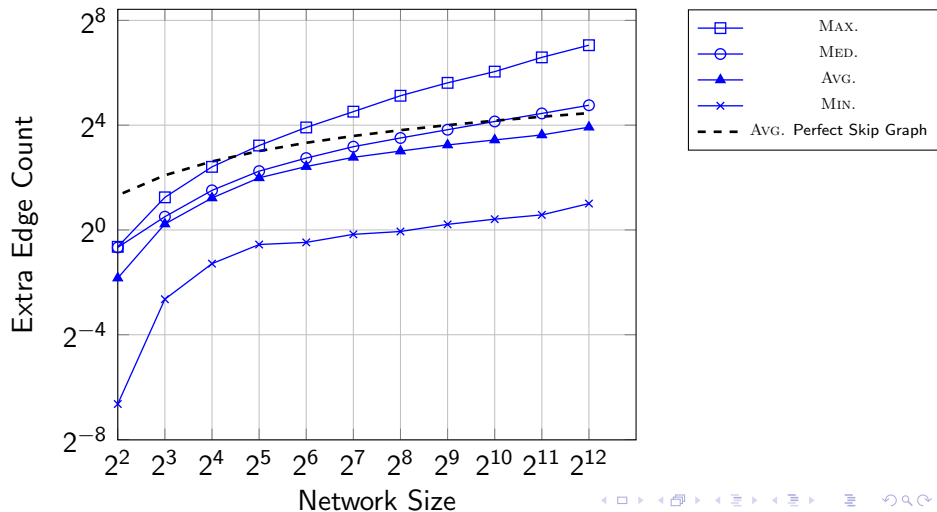
# MULTISKIPGRAPH

The self-stabilizing part: updating local variables

- $LeftLevel[i]$ : the left level- $i$  neighbor
- $RightLevel[i]$ : the right level- $i$  neighbor
- $LeftUnknown$ : the set of left neighbors which are not assigned to any level, marked as *unknown*.
- $RightUnknown$ : the set of right neighbors which are not assigned to any level, marked as *unknown*.
- $Left$ : the set of all left neighbors
- $Right$ : the set of all right neighbors

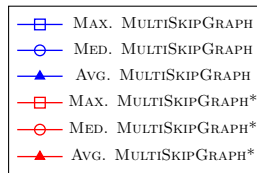
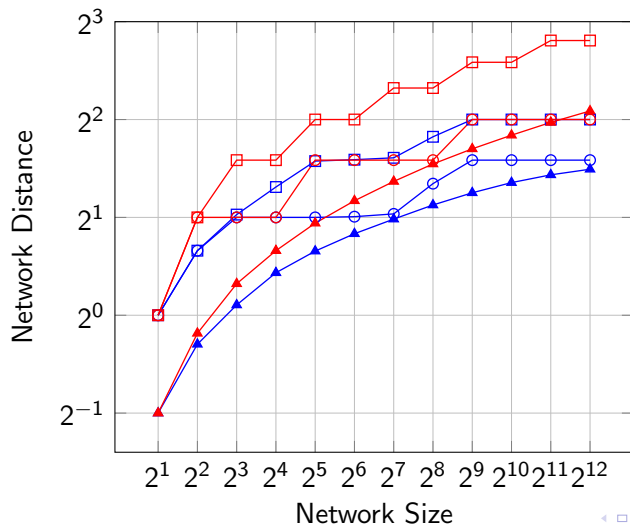
# Comparison in Stabilization

## Extra Edge Count



# Comparison in Stabilization

Network distance between two nodes



# Comparison in Searchability

## Hop Count

