

# ElastO: Efficient Maintenance of Scalable Overlays for Topic-based Pub/Sub under Churn [Submission to Middleware '14]

**Chen Chen**

Joint work with Roman Vitenberg and Hans-Arno Jacobsen

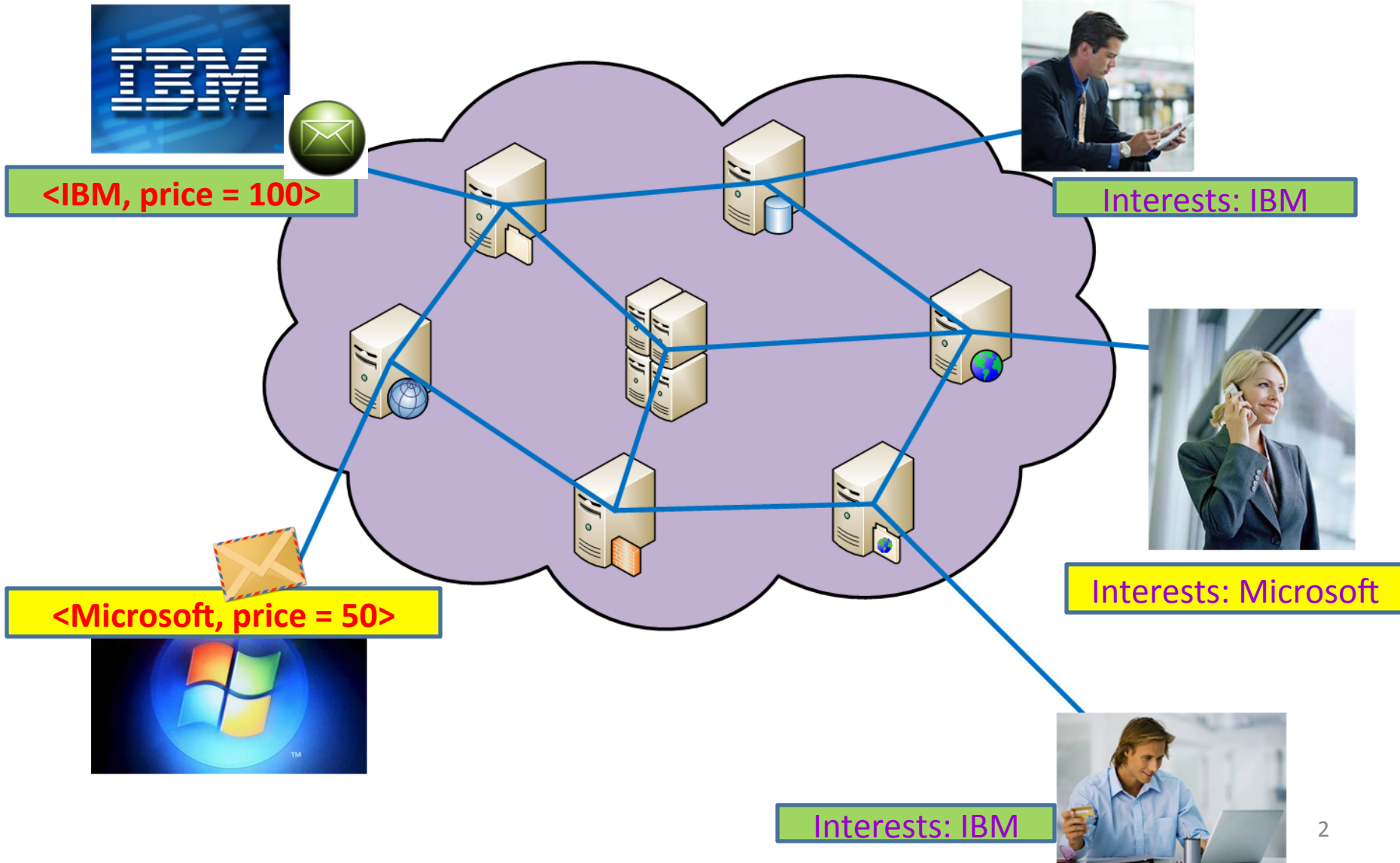


The Edward S. Rogers Sr. Department  
of Electrical & Computer Engineering  
**UNIVERSITY OF TORONTO**

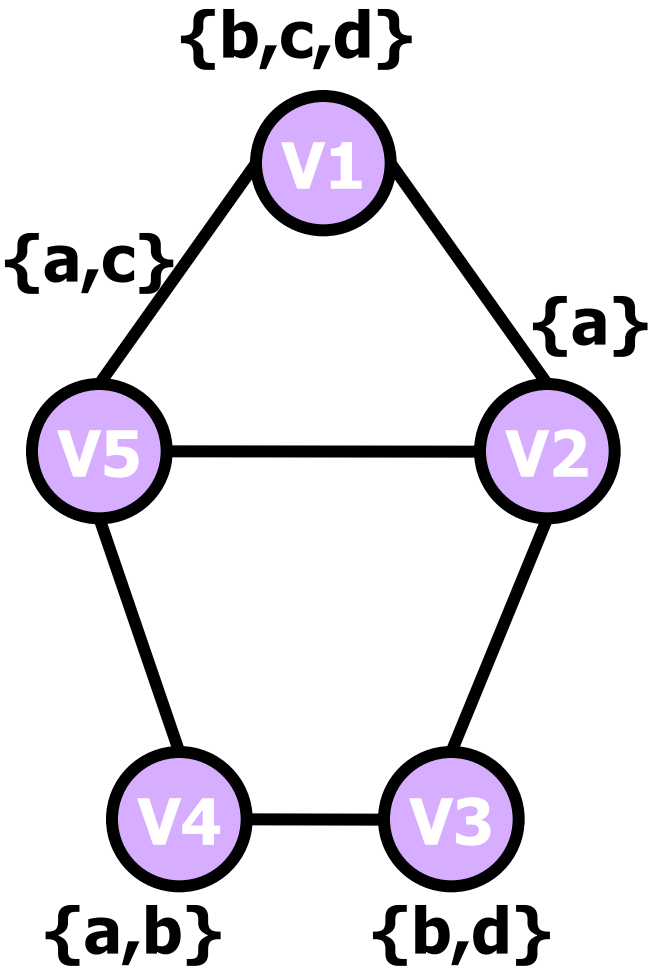


**UNIVERSITY  
OF OSLO**

# Publish/Subscribe (pub/sub)

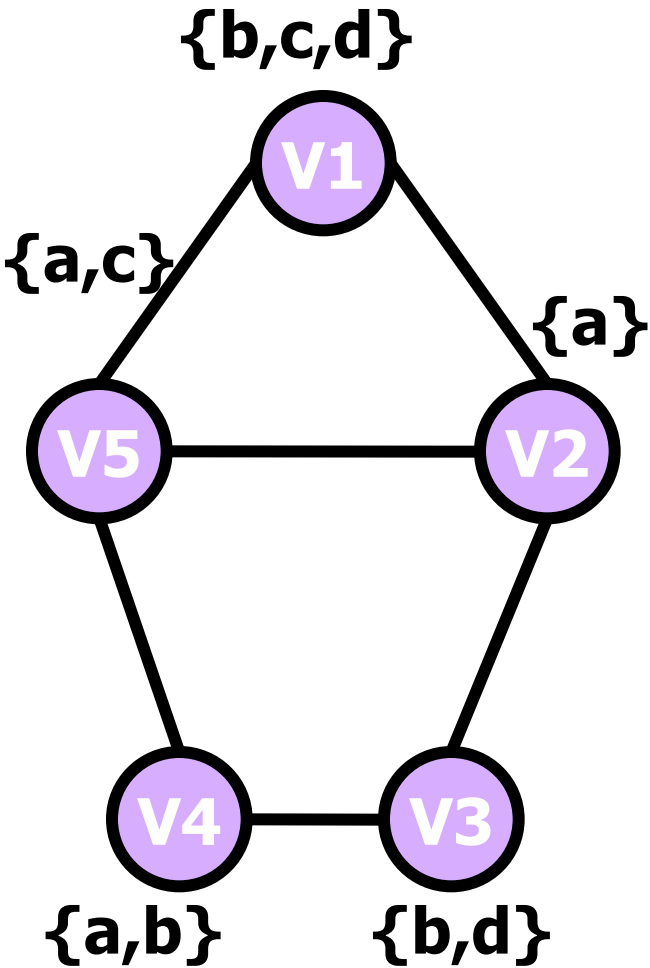


# Topic-connected overlay (TCO)

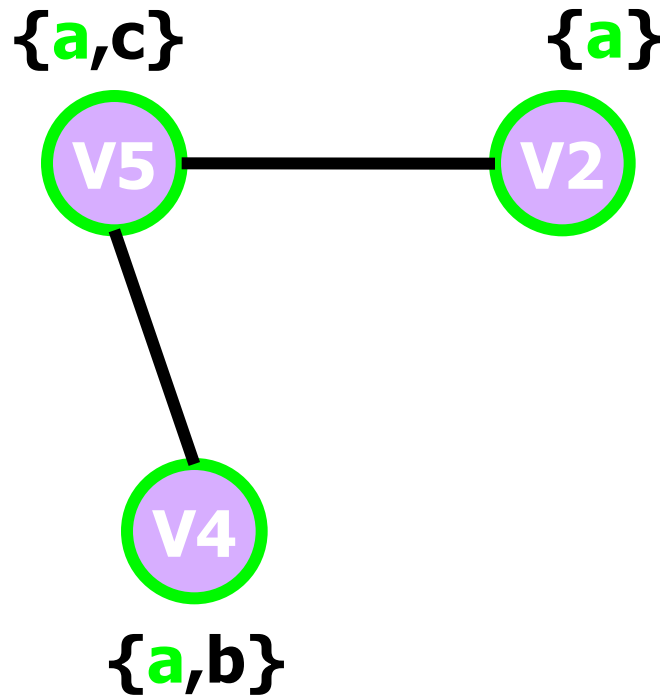


An overlay G

# Topic-connected overlay (TCO)

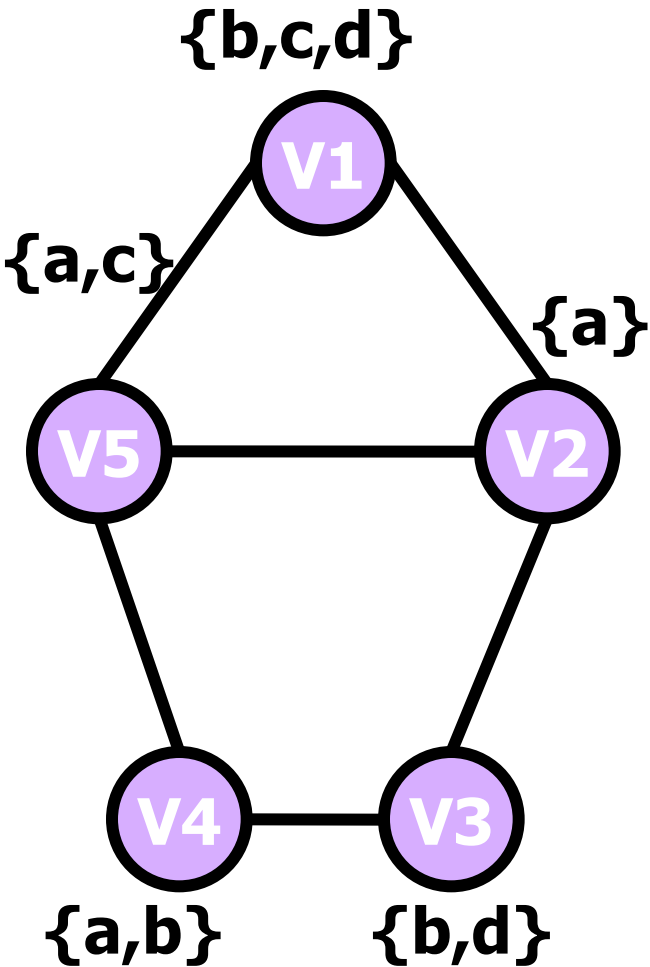


An overlay  $G$

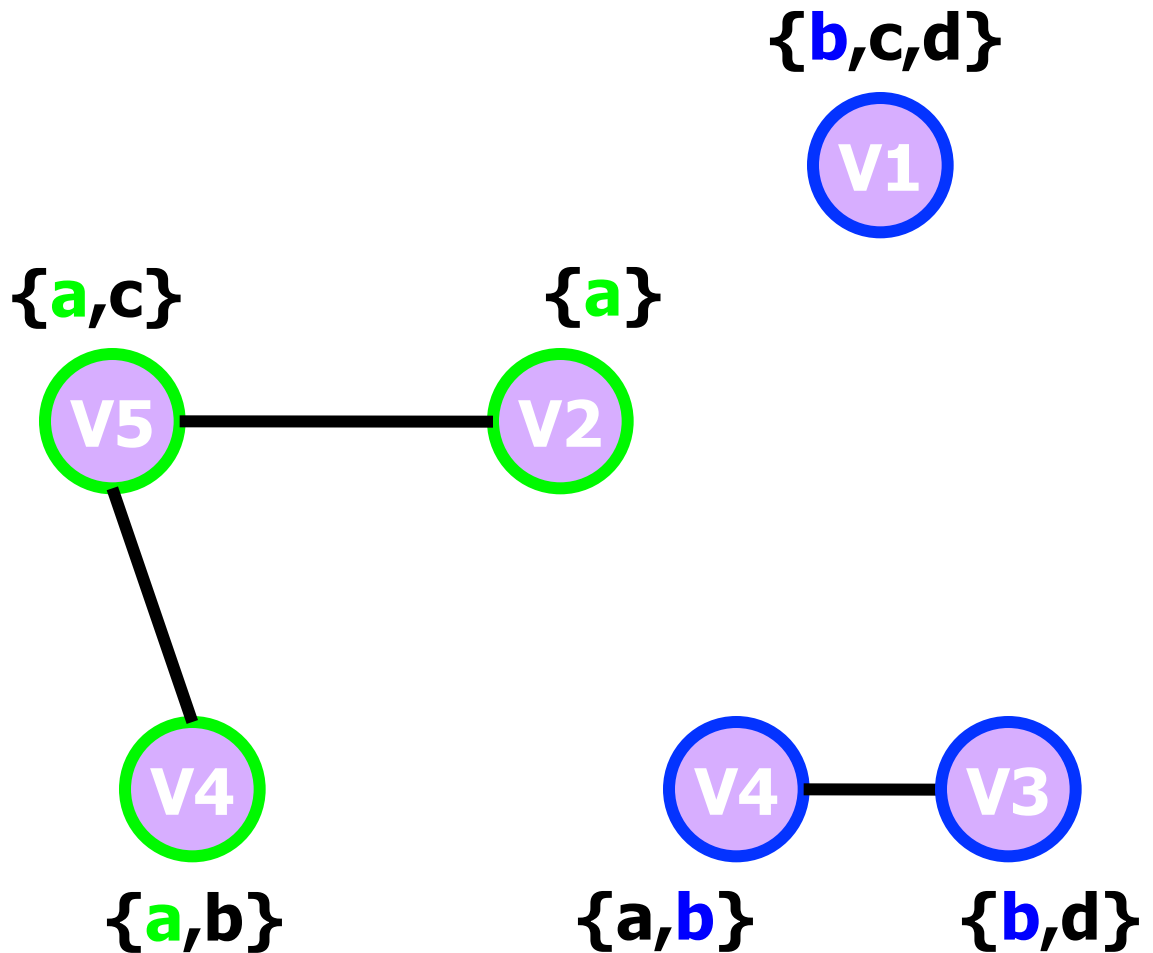


$G_a$  is topic-connected with one *TC-component*

# Topic-connected overlay (TCO)



An overlay  $G$



$G_a$  is topic-connected with one *TC-component*

$G_b$  is NOT topic-connected with two *TC-components*

# Approaches to build pub/sub TCO

		Knowledge	Churn Handling	Time	Average degree	Maximum degree
<b>Centralized algorithms</b>	Low-ODA	<b>Global</b>	<b>X</b>	<b>Slow</b>	$O(\rho \ln  V  T )$	$O( V /\rho \ln  V  T )$
	GM				$O(\ln  V  T )$	$\Theta( V )$
	MinMax-ODA				$\Theta( V )$	$O(\ln  V  T )$
	DC				$O(p \ln  V  T )$	$\Theta( V )$
	DCBR-M				$\Theta( V )$	$O(\eta + \ln  V  T )$
	GM <sub>2</sub>				$O(U + \ln  V  T )$	$\Theta( V )$
<b>Decentralized protocols</b>	SpiderCast, PolderCast, StAN, etc.	<b>Local/Global</b>	<b>✓</b>	<b>Fast</b>	<b>Unknown</b>	<b>Unknown</b>

# Our hybrid solution: ElastO

		Knowledge	Churn Handling	Time	Average degree	Maximum degree
<b>Hybrid</b>	<b>ElastO</b>	<b>Local</b>	✓	<b>Fast</b>	$\approx O(\rho \ln  V  T )$	$\approx O( V /\rho \ln  V  T )$
<b>Centralized algorithms</b>	Low-ODA	<b>Global</b>	<b>✗</b>	<b>Slow</b>	$O(\rho \ln  V  T )$	$O( V /\rho \ln  V  T )$
	GM				$O(\ln  V  T )$	$\Theta( V )$
	MinMax-ODA				$\Theta( V )$	$O(\ln  V  T )$
	DC				$O(p \ln  V  T )$	$\Theta( V )$
	DCBR-M				$\Theta( V )$	$O(\eta + \ln  V  T )$
	GM2				$O(U + \ln  V  T )$	$\Theta( V )$
<b>Decentralized protocols</b>	SpiderCast, PolderCast, StAN, etc.	<b>Local/Global</b>	✓	<b>Fast</b>	<b>Unknown</b>	<b>Unknown</b>

# Requirements

## A. **Overlay quality**

- TCO, low node degrees, small diameters
- Close to centralized algorithms

## B. **Responsiveness**

- Comparable to decentralized protocols

## C. **Scalability of the decentralized solution**

- No centralized control
- Partial and local view
- Churn handling only impacts a small portion of nodes

## D. **Fairness and load balancing**

- Computation, communication, storage

## E. **Reliability** against concurrent churn events



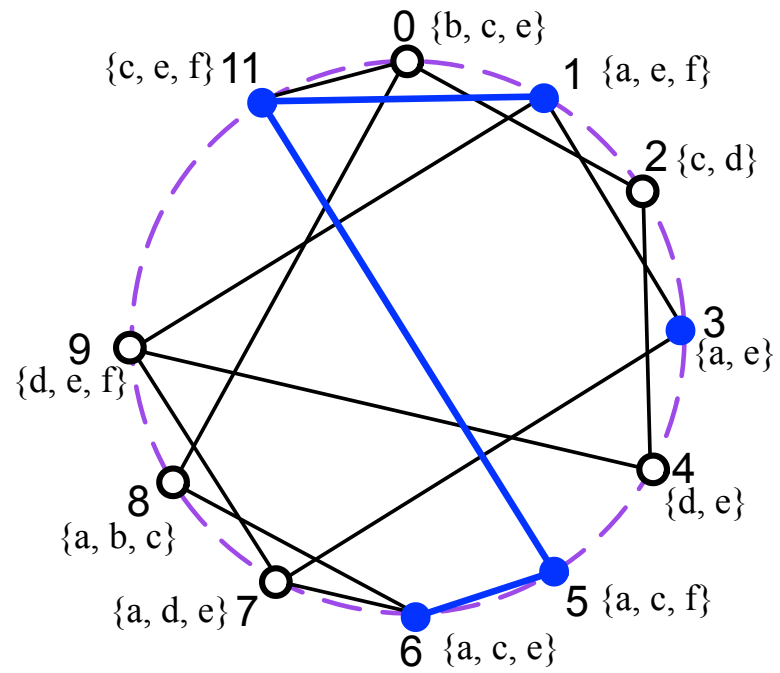
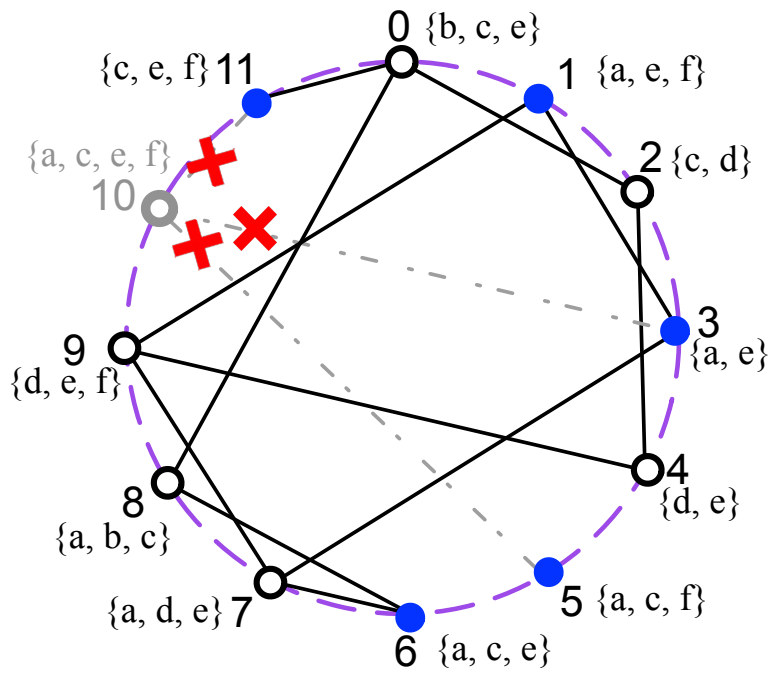
# Our contributions - ElastO

- Complete architecture and protocol design
- **Shadow**-based strategies
  - TCO recovery under churn
- **Primary-backup** mechanism
  - support shadow sets in decentralized systems
- Gossip-based peer sampling service with unique **local view selection** algorithms
  - build and deploy backup sets for all nodes
- Comprehensive evaluation

# Shadow-based strategies

Upon each churn event

- Select a *proper* subset of nodes, namely the **shadow set**, for TCO recovery



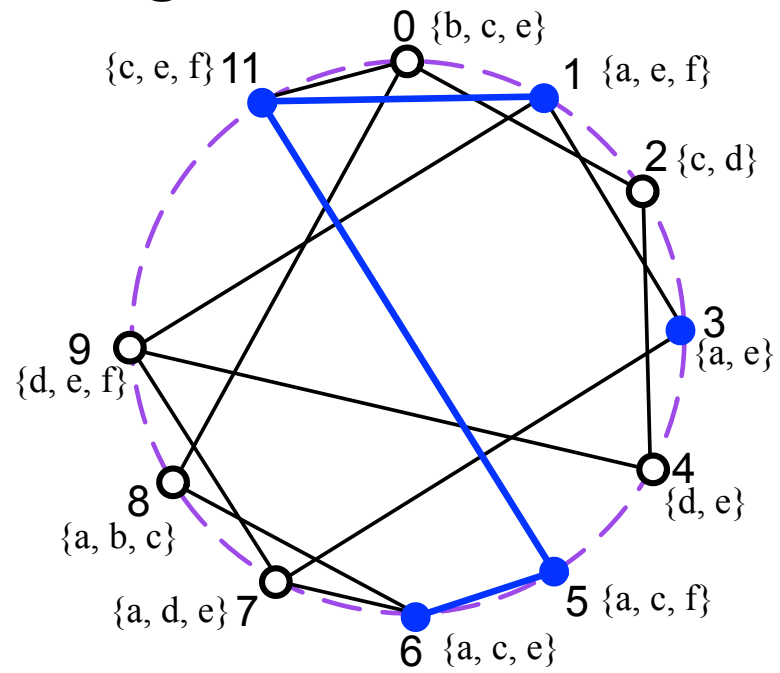
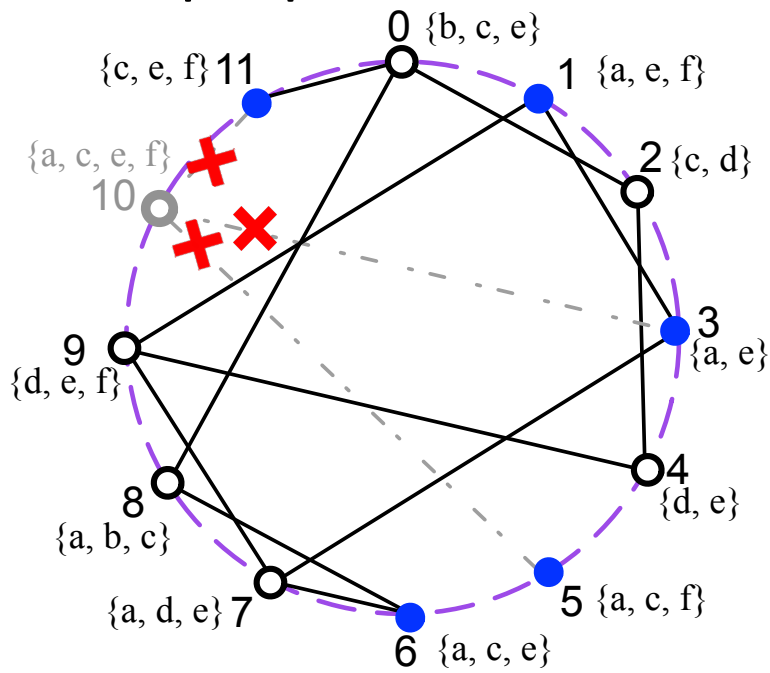
# Pre-compute the shadow set primary-backup mechanism

- ElastO needs to remember

- 1) TCO neighbors,  $v_{10} \cdot \mathcal{N} = \{v_3, v_5, v_{11}\}$

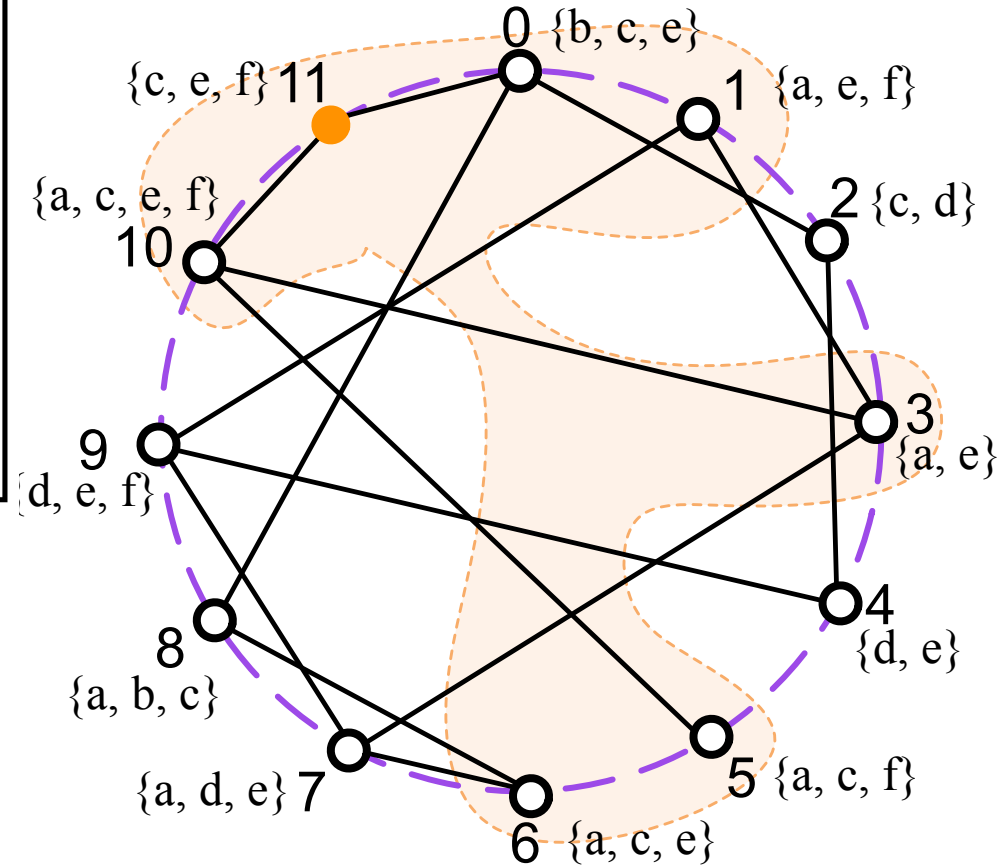
- 2) backup set,  $B(v_{10}) = \{v_1, v_6\}$

in preparation of  $v_{10}$ 's leaving



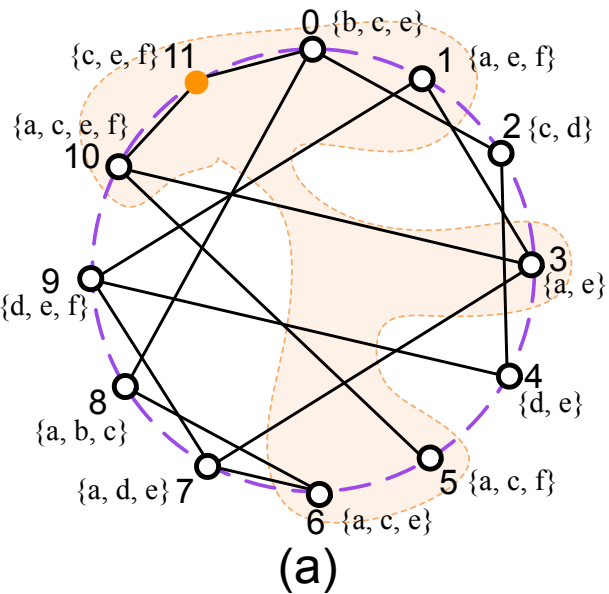
# Local view of each node

$.self$	$v_{11}$
$.\mathcal{N}$	$\{v_0, v_{10}\}$
$\mathcal{D}$	$\emptyset$
$.succ[1]$	$v_0$
$.pred[1]$	$v_{10}$
$.pneighb[1]$	$\{v_3, v_5, v_{11}\}$
$.backup[1]$	$\{v_1, v_6\}$



# Example: initially stable

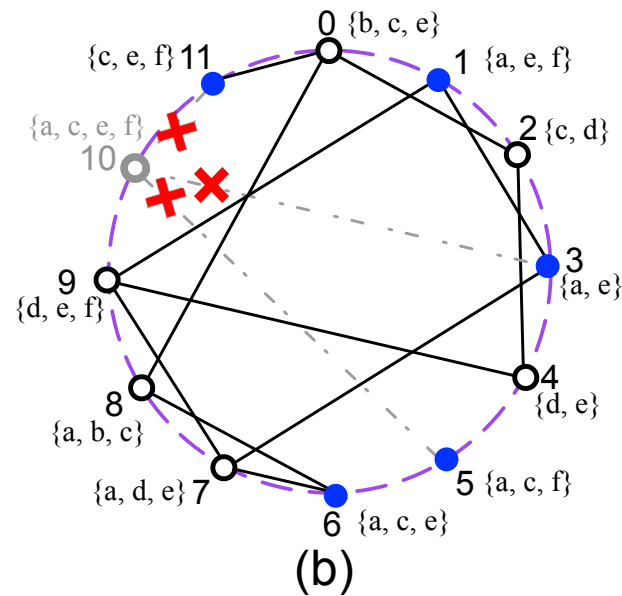
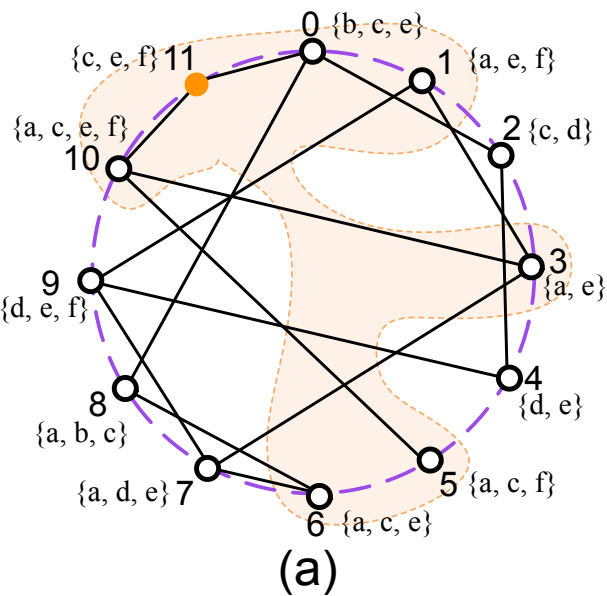
<i>.self</i>	$v_{11}$
<i>.N</i>	$\{v_0, v_{10}\}$
<i>.D</i>	$\emptyset$
<i>.succ</i> [1]	$v_0$
<i>.pred</i> [1]	$v_{10}$
<i>.pneighb</i> [1]	$\{v_3, v_5, v_{11}\}$
<i>.backup</i> [1]	$\{v_1, v_6\}$



# Example: node departure

<i>.self</i>	$v_{11}$
<i>.N</i>	$\{v_0, v_{10}\}$
<i>.D</i>	$\emptyset$
<i>.succ</i> [1]	$v_0$
<i>.pred</i> [1]	$v_{10}$
<i>.pneighb</i> [1]	$\{v_3, v_5, v_{11}\}$
<i>.backup</i> [1]	$\{v_1, v_6\}$

<i>.self</i>	$v_{11}$
<i>.N</i>	$\{v_0, v_{10}\}$
<i>.D</i>	$\{v_{10}\}$
<i>.succ</i> [1]	$v_0$
<i>.pred</i> [1]	$v_{10}$
<i>.pneighb</i> [1]	$\{v_3, v_5, v_{11}\}$
<i>.backup</i> [1]	$\{v_1, v_6\}$

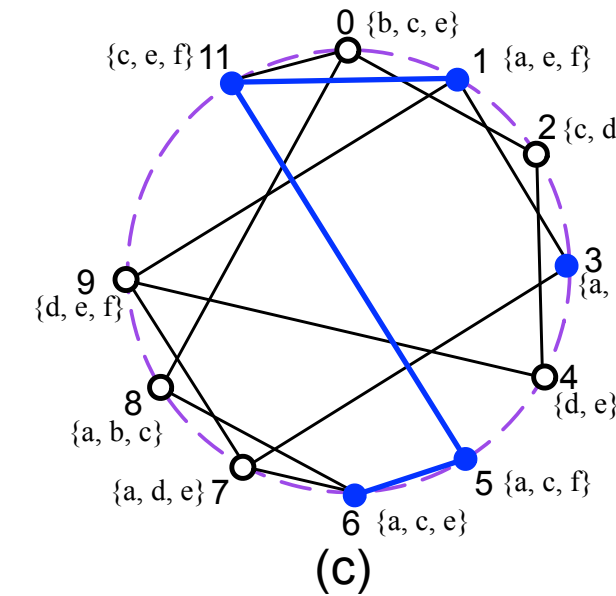
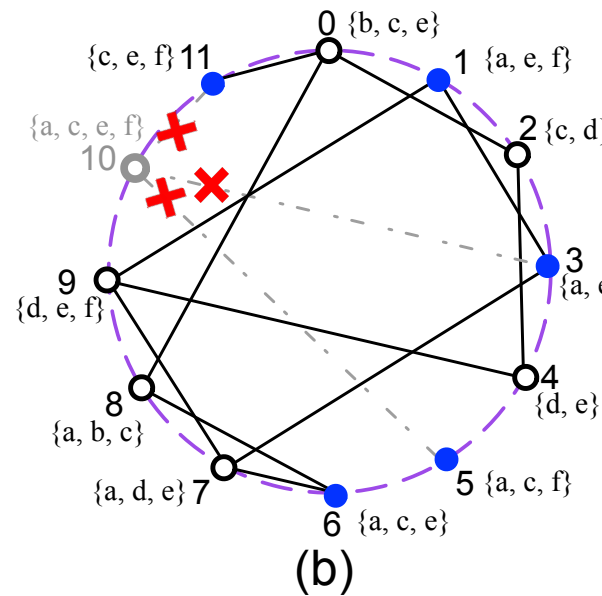
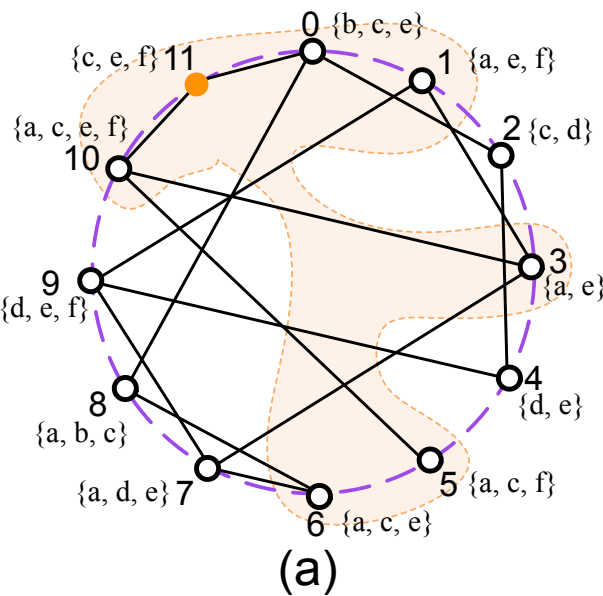


# Example: node departure handled

<i>.self</i>	$v_{11}$
<i>.N</i>	$\{v_0, v_{10}\}$
<i>.D</i>	$\emptyset$
<i>.succ</i> [1]	$v_0$
<i>.pred</i> [1]	$v_{10}$
<i>.pneighb</i> [1]	$\{v_3, v_5, v_{11}\}$
<i>.backup</i> [1]	$\{v_1, v_6\}$

<i>.self</i>	$v_{11}$
<i>.N</i>	$\{v_0, v_{10}\}$
<i>.D</i>	$\{v_{10}\}$
<i>.succ</i> [1]	$v_0$
<i>.pred</i> [1]	$v_{10}$
<i>.pneighb</i> [1]	$\{v_3, v_5, v_{11}\}$
<i>.backup</i> [1]	$\{v_1, v_6\}$

<i>.self</i>	$v_{11}$
<i>.N</i>	$\{v_0, v_1, v_5\}$
<i>.D</i>	$\{\}$
<i>.succ</i> [1]	$v_0$
<i>.pred</i> [1]	$v_9$
<i>.pneighb</i> [1]	$\{v_1, v_4, v_7\}$
<i>.backup</i> [1]	$\{v_5, v_7\}$



# Evaluation: experimental setup

- Real-world pub/sub workloads

facebook

twitter 

- Synthetic pub/sub workloads

Expo, Zipf, Unif

- Churn traces

Google



# Evaluation: algorithms and protocols

---

**ElastO**                      **Our proposed system**

ElastO-L                      Local view at each node

ElastO-G                      Global view at each node

---

**LowODA**                      **Low Max and Avg Degree Overlay Design Algorithm**

LowODA-Inc                      Incrementally repair TCO regarding existing links

LowODA-Re                      Reconstruct TCO from scratch regardless of existing links

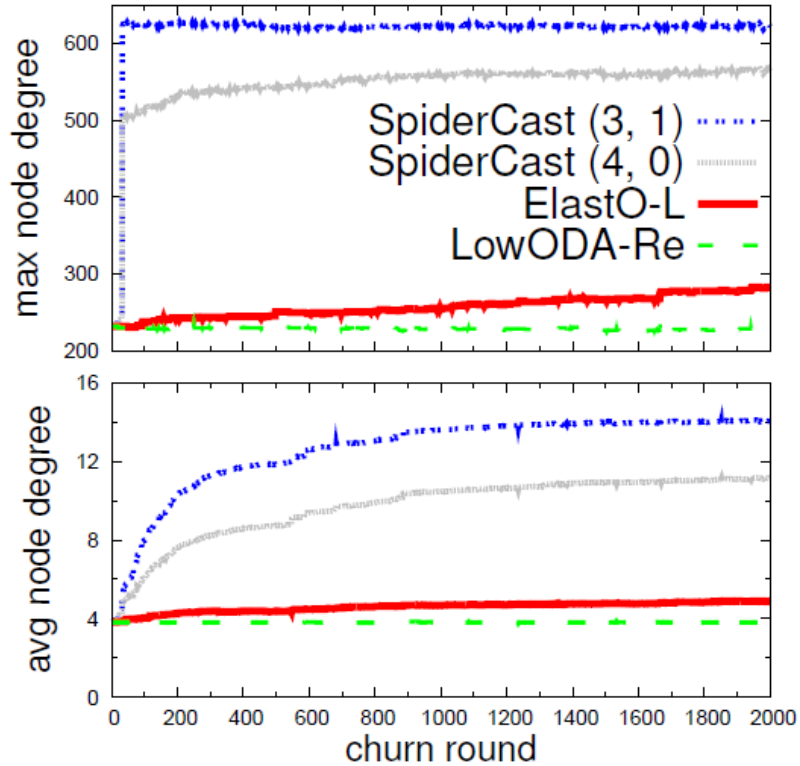
---

**SpiderCast**                      **A peer-to-peer protocol to build TCO**

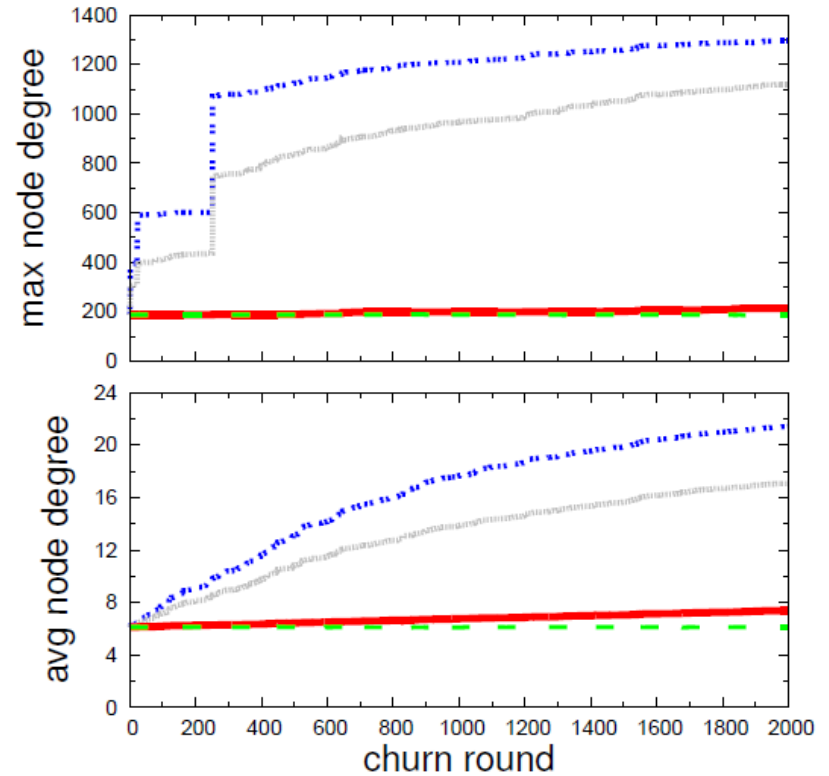
SpiderCast(Kg,Kr)                      Two neighbor selection heuristics: greedy and random

---

# Evaluation: node degrees under churn

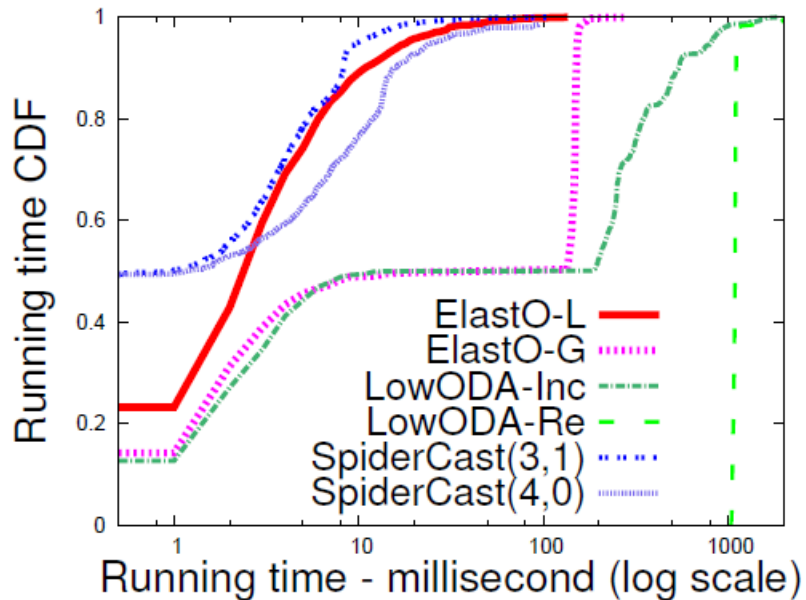


(a) FB 1K

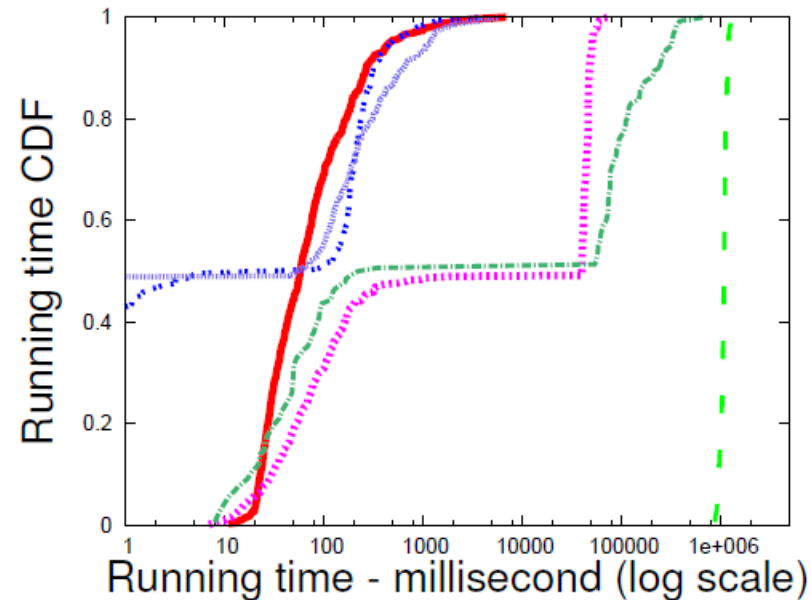


(b) FB 10K

# Evaluation: runtime cost



(c) TW 1K



(d) TW 10K

# Conclusion

		Knowledge	Churn Handling	Time	Average degree	Maximum degree
<b>Hybrid</b>	<b>ElastO</b>	<b>Local</b>	✓	<b>Fast</b>	$\approx O(\rho \ln  V   T )$	$\approx O( V /\rho \ln  V   T )$
<b>Centralized algorithms</b>	Low-ODA	<b>Global</b>	<b>✗</b>	<b>Slow</b>	$O(\rho \ln  V   T )$	$O( V /\rho \ln  V   T )$
	GM				$O(\ln  V   T )$	$\Theta( V )$
	MinMax-ODA				$\Theta( V )$	$O(\ln  V   T )$
	<b>DC</b>				$O(p \ln  V   T )$	$\Theta( V )$
	<b>DCBR-M</b>				$\Theta( V )$	$O(\eta + \ln  V   T )$
	<b>GM<sub>2</sub></b>				$O(U + \ln  V   T )$	$\Theta( V )$
<b>Decentralized protocols</b>	SpiderCast, PolderCast, etc.	<b>Local/Global</b>	✓	<b>Fast</b>	<b>Unknown</b>	<b>Unknown</b>

# Current research in IBM

## Distributed pub/sub for federated messaging and IoT

- Scalability: 100 Million users/topics/subscriptions
- Quality of Service (QoS) for MQTT
  - 0: at most once
  - 1: at least once
  - 2: exactly once
- Topics and subscriptions
  - hierarchy and wildcards
  - e.g., sensors/+/temperature/+
- Overlay topology and routing protocols
- Membership management