

FAILURE CONTROL IN MULTIPATH ROUTE TRACING

Darryl Veitch¹

dveitch@unimelb.edu.au

.....

Collaborators

Brice Augustin² Renata Teixeira² Timur Friedman²

¹CUBIN, Department of Electrical & Electronic Engineering, University of Melbourne, Australia

²UPMC Paris Universit as, and CNRS, LIP6 Laboratory, Paris, France

TRACEROUTE, A USEFUL TOOL

THE TOOL

- Aims to discover IP-level path between source and destination
- Time to Live (TTL) field in IP header targets hop-depth
- Returning TTL expiry packets carry address of last interface

ITS APPLICATIONS

- Diagnosing network problems
- Drawing network maps

TRACEROUTE, A USEFUL TOOL

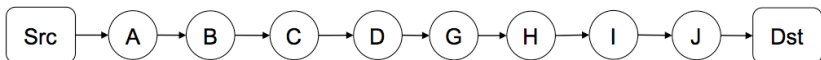
THE TOOL

- Aims to discover IP-level path between source and destination
- Time to Live (TTL) field in IP header targets hop-depth
- Returning TTL expiry packets carry address of last interface

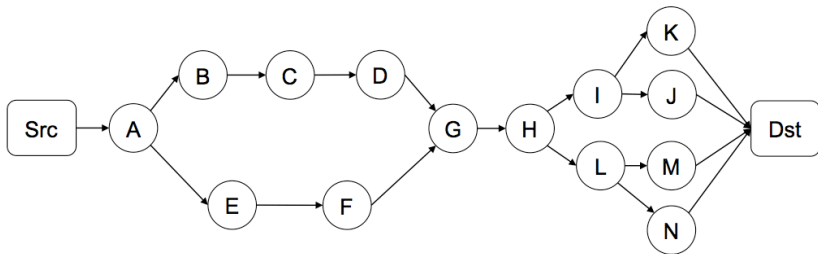
ITS APPLICATIONS

- Diagnosing network problems
- Drawing network maps

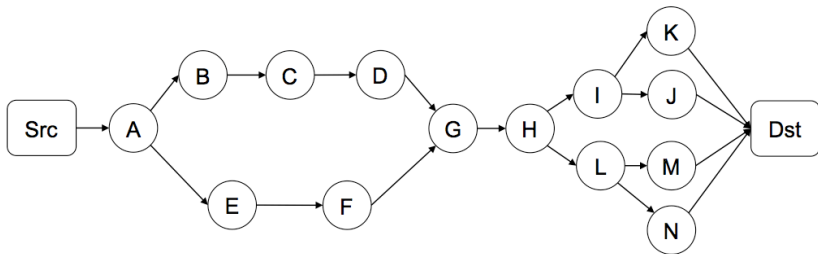
CLASSIC TRACEROUTE IGNORES LOAD BALANCING



CLASSIC TRACEROUTE IGNORES LOAD BALANCING



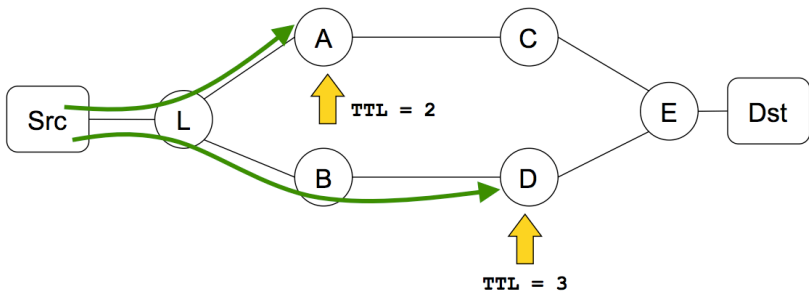
CLASSIC TRACEROUTE IGNORES LOAD BALANCING



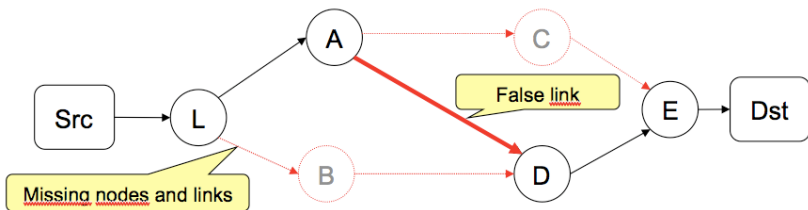
IMPACT ON CLASSIC TRACEROUTE

- Paths missed
- Inferred paths erroneous

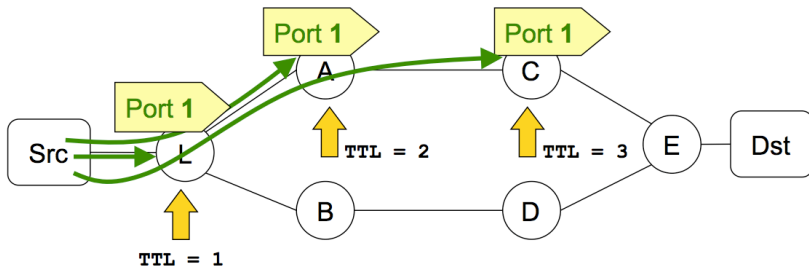
CLASSIC FAILURES



CLASSIC FAILURES



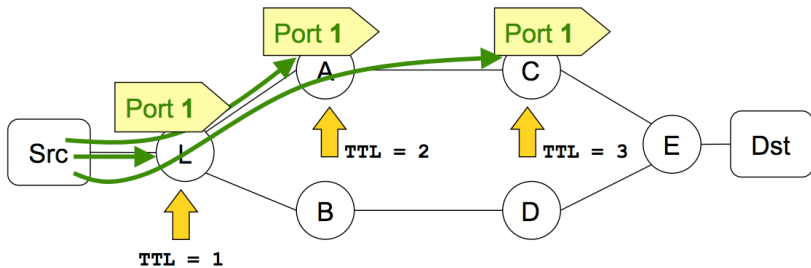
PARIS TRACEROUTE SOLUTION FOR PER-FLOW LB



EXPLICITLY CONTROL FLOW KEYS

- Ensures probes follow a fixed path assuming per-flow LB
- Failure to do so reveals presence of per-packet LB

PARIS TRACEROUTE SOLUTION FOR PER-FLOW LB



EXPLICITLY CONTROL FLOW KEYS

- Ensures probes follow a fixed path assuming per-flow LB
- Failure to do so reveals presence of per-packet LB

PARIS TRACEROUTE DEVELOPMENT

OUR PRIOR WORK

- **IMC'06:**
 - Revealing flaws in Classic Traceroute
 - Introduction of flow key fixing idea for single path under per-flow LB
- **E2EMon'07:** Introduction of
 - Stochastic probing algorithm **Multipath detection algorithm** (MDA)
 - Per-node significance levels
- **IMC'07:** Use of MDA to study LB prevalence in Internet

THIS PAPER

- Multipath **failure** probabilities, calculation and bounds
- Extension from node to path level statistical guarantees
- Exact significance levels for multipath discovery

PARIS TRACEROUTE DEVELOPMENT

OUR PRIOR WORK

- **IMC'06:**
 - Revealing flaws in Classic Traceroute
 - Introduction of flow key fixing idea for single path under per-flow LB
- **E2EMon'07:** Introduction of
 - Stochastic probing algorithm **Multipath detection algorithm** (MDA)
 - Per-node significance levels
- **IMC'07:** Use of MDA to study LB prevalence in Internet

THIS PAPER

- Multipath **failure** probabilities, calculation and bounds
- Extension from node to path level statistical guarantees
- Exact significance levels for multipath discovery

MULTIPATH ROUTE DISCOVERY

ADAPTIVE PROBING STRATEGY: AT EACH HOP

- Controlled exploration of multipath by varying flow identifier of probes
- Send enough probes to enumerate all interfaces with high confidence (retransmit if necessary)
- Classify load balancers: **per-flow** or per-packet

UNDERSTANDING FAILURE

FAILURE

- Means not finding **all** links
- Operates at node and graph level
- **Implies topology known!** but we want to *discover* multipath...
- Importance of failure:
 - gives control for fixed (common) topologies
 - provides the foundation for everything

ASSUMPTIONS FOR CALCULATION

- Node level: probe returns interface IID uniformly out of $\{1, 2, \dots, K\}$
- Graph level: all nodes independent

UNDERSTANDING FAILURE

FAILURE

- Means not finding **all** links
- Operates at node and graph level
- **Implies topology known!** but we want to *discover* multipath...
- Importance of failure:
 - gives control for fixed (common) topologies
 - provides the foundation for everything

ASSUMPTIONS FOR CALCULATION

- Node level: probe returns interface IID uniformly out of $\{1, 2, \dots, K\}$
- Graph level: all nodes independent

UNDERSTANDING FAILURE

FAILURE

- Means not finding **all** links
- Operates at node and graph level
- **Implies topology known!** but we want to *discover* multipath...
- Importance of failure:
 - gives control for fixed (common) topologies
 - provides the foundation for everything

ASSUMPTIONS FOR CALCULATION

- Node level: probe returns interface IID uniformly out of $\{1, 2, \dots, K\}$
- Graph level: all nodes independent

FAILURE PROBABILITIES

NODE FAILURE PROBABILITY β_K

Define stopping points $\{n_k\}$ for each of $k \in \{1, 2, \dots, K\}$

Corresponding *stopping probability* t_K^k at level k ($\sum_{k=1}^K t_K^k = 1$)

$$\beta_1 = 0$$

$$\beta_K = \sum_{k=1}^{K-1} t_K^k, \quad K > 1$$

FAILURE PROBABILITIES

NODE FAILURE PROBABILITY β_K

Define stopping points $\{n_k\}$ for each of $k \in \{1, 2, \dots, K\}$

Corresponding *stopping probability* t_K^k at level k ($\sum_{k=1}^K t_K^k = 1$)

$$\beta_1 = 0$$

$$\beta_K = \sum_{k=1}^{K-1} t_K^k, \quad K > 1$$

FAILURE PROBABILITIES

NODE FAILURE PROBABILITY β_K

Define stopping points $\{n_k\}$ for each of $k \in \{1, 2, \dots, K\}$

Corresponding *stopping probability* t_K^k at level k ($\sum_{k=1}^K t_K^k = 1$)

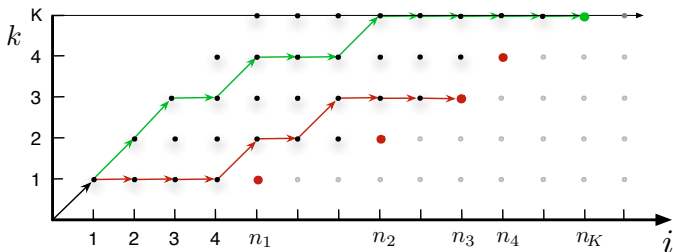
$$\beta_1 = 0$$

$$\beta_K = \sum_{k=1}^{K-1} t_K^k, \quad K > 1$$

GRAPH FAILURE PROBABILITY β_{all}

$$\beta_{\text{all}} = 1 - \prod_i (1 - \beta_{K_i})$$

NODE STATE SPACE AND TRAJECTORIES



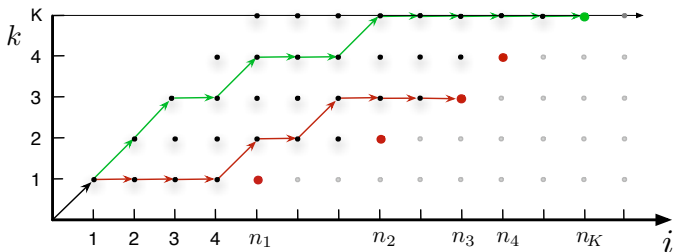
- Transition probabilities after new probe:

– Find new: $q_k = (K - k)/K$

– Nothing new: $p_k = 1 - q_k = k/K$

- To reach $(n_1, 1)$: $t_K^1 = p_1^{n_1-1}$; For $(n_2, 2)$: $t_K^2 = \sum_{j_1=0}^{n_1-2} p_1^{j_1} q_1 p_2^{n_2-2-j_1}$

NODE STATE SPACE AND TRAJECTORIES



- Transition probabilities after new probe:
 - Find new: $q_k = (K - k)/K$
 - Nothing new: $p_k = 1 - q_k = k/K$
- To reach $(n_1, 1)$: $t_K^1 = p_1^{n_1-1}$; For $(n_2, 2)$: $t_K^2 = \sum_{j_1=0}^{n_1-2} p_1^{j_1} q_1 p_2^{n_2-2-j_1}$

CONDITIONAL FAILURE PROBABILITIES

Stopping probabilities t_K^k are unconditional.

For **adaptive** algorithm, want **conditional** ones *given j interfaces already seen*.

$$\begin{aligned} t_{K,j}^k &= \Pr(\text{stop at level } k \mid \text{didn't stop at levels below } j) \\ &= \frac{\Pr(\text{stop at level } k)}{\Pr(\text{didn't stop at levels below } j)} = \frac{t_K^k}{\sum_{l=j}^K t_K^l} \end{aligned}$$

for $k \geq j$ and zero otherwise.

Corresponding conditional failure probability, $1 \leq j \leq K$, is

$$\begin{aligned} \beta_{K,K} &= 0 \\ \beta_{K,j} &= \sum_{k=1}^{K-1} t_{K,j}^k = \sum_{k=j}^{K-1} t_{K,j}^k = \frac{\sum_{k=j}^{K-1} t_K^k}{\sum_{l=j}^K t_K^l}, \quad j < K, \end{aligned}$$

and naturally $\beta_{K,1} = \beta_K$.

CONDITIONAL FAILURE PROBABILITIES

Stopping probabilities t_K^k are unconditional.

For **adaptive** algorithm, want **conditional** ones *given j interfaces already seen*.

$$\begin{aligned} t_{K,j}^k &= \Pr(\text{stop at level } k \mid \text{didn't stop at levels below } j) \\ &= \frac{\Pr(\text{stop at level } k)}{\Pr(\text{didn't stop at levels below } j)} = \frac{t_K^k}{\sum_{l=j}^K t_K^l} \end{aligned}$$

for $k \geq j$ and zero otherwise.

Corresponding conditional failure probability, $1 \leq j \leq K$, is

$$\begin{aligned} \beta_{K,K} &= 0 \\ \beta_{K,j} &= \sum_{k=1}^{K-1} t_{K,j}^k = \sum_{k=j}^{K-1} t_{K,j}^k = \frac{\sum_{k=j}^{K-1} t_K^k}{\sum_{l=j}^K t_K^l}, \quad j < K, \end{aligned}$$

and naturally $\beta_{K,1} = \beta_K$.

CONDITIONAL FAILURE PROBABILITIES

Stopping probabilities t_K^k are unconditional.

For **adaptive** algorithm, want **conditional** ones *given j interfaces already seen*.

$$\begin{aligned} t_{K,j}^k &= \Pr(\text{stop at level } k \mid \text{didn't stop at levels below } j) \\ &= \frac{\Pr(\text{stop at level } k)}{\Pr(\text{didn't stop at levels below } j)} = \frac{t_K^k}{\sum_{l=j}^K t_K^l} \end{aligned}$$

for $k \geq j$ and zero otherwise.

Corresponding conditional failure probability, $1 \leq j \leq K$, is

$$\begin{aligned} \beta_{K,K} &= 0 \\ \beta_{K,j} &= \sum_{k=1}^{K-1} t_{K,j}^k = \sum_{k=j}^{K-1} t_{K,j}^k = \frac{\sum_{k=j}^{K-1} t_K^k}{\sum_{l=j}^K t_K^l}, \quad j < K, \end{aligned}$$

and naturally $\beta_{K,1} = \beta_K$.

HYPOTHESIS TESTING FRAMEWORK

In practice, K unknown. Role of algorithm is to give up search reluctantly:

Control probability α_k of stopping at state (n_k, k) , given have k interfaces already, under null hypothesis that in fact are $K' > k$ of them.

Null-hyp is **composite**, but since $p_{k,K'} = k/K'$, sig-level set by $K' = k + 1$.

SIMPLE NULL HYPOTHESIS OF $K' = k + 1$

- Stopping at (n_k, k) and failing to find new successors is a type-I error
- Stopping accepts alternate hypothesis that $K' = k$
- Test has optimal power! (if $k = K' = K$, impossible to find more)
- Significance level is

$$\alpha_k = t_{k+1,k}^k = \frac{t_{k+1}^k}{t_{k+1}^k + t_{k+1}^{k+1}}, \quad 1 \leq k \leq K.$$

HYPOTHESIS TESTING FRAMEWORK

In practice, K unknown. Role of algorithm is to give up search reluctantly:

Control probability α_k of stopping at state (n_k, k) , given have k interfaces already, under null hypothesis that in fact are $K' > k$ of them.

Null-hyp is **composite**, but since $p_{k,K'} = k/K'$, sig-level set by $K' = k + 1$.

SIMPLE NULL HYPOTHESIS OF $K' = k + 1$

- Stopping at (n_k, k) and failing to find new successors is a type-I error
- Stopping accepts alternate hypothesis that $K' = k$
- Test has optimal power! (if $k = K' = K$, impossible to find more)
- Significance level is

$$\alpha_k = t_{k+1,k}^k = \frac{t_{k+1}^k}{t_{k+1}^k + t_{k+1}^{k+1}}, \quad 1 \leq k \leq K.$$

CONNECTING SIGNIFICANCE AND FAILURE

Upon terminating a node-level discovery, a significance level will be reported.
How does it relate to the true (conditional) probability of error?

RELATIONSHIP NOT SIMPLE!

- If $k = K$ found, then $\alpha_K > \beta_{K,K} = 0$, so algorithm is conservative
- If $k = K - 1$, algo uses $K^* = K$, and $\alpha_{K-1} = \beta_{K,K-1}^* = \beta_{K,K}$, ideal!
- If $k < K - 1$, can go either way, e.g., with $K = 3$ and $k = 1$:

• If $\beta_{3,1} > \beta_{3,2}$, then $\alpha_1 > \alpha_2$, so algorithm is conservative

• If $\beta_{3,1} < \beta_{3,2}$, then $\alpha_1 < \alpha_2$, so algorithm is liberal

CONNECTING SIGNIFICANCE AND FAILURE

Upon terminating a node-level discovery, a significance level will be reported.
How does it relate to the true (conditional) probability of error?

RELATIONSHIP NOT SIMPLE!

- If $k = K$ found, then $\alpha_K > \beta_{K,K} = 0$, so algorithm is conservative
- If $k = K - 1$, algo uses $K' = K$, and $\alpha_{K-1} = \beta_{K,K-1}^{K-1} = \beta_{K,K-1}$, ideal!
- If $k < K - 1$, can go either way, e.g. with $K = 3$ and $k = 1$:
 - If $(n_1, n_2, n_3) = (2, 30, 30)$ then $\alpha_1 > \beta_{3,1}$ (pessimistic, conservative)
 - If $(n_1, n_2, n_3) = (2, 3, 4)$, find $\alpha_1 < \beta_{3,1}$ (optimistic, **misleading!**)

CONNECTING SIGNIFICANCE AND FAILURE

Upon terminating a node-level discovery, a significance level will be reported.
How does it relate to the true (conditional) probability of error?

RELATIONSHIP NOT SIMPLE!

- If $k = K$ found, then $\alpha_K > \beta_{K,K} = 0$, so algorithm is conservative
- If $k = K - 1$, also uses $K' = K$, and $\alpha_{K-1} = t_{K,K-1}^{K-1} = \beta_{K,K-1}$, ideal!
- If $k < K - 1$, can go either way, e.g. with $K = 3$ and $k = 1$:
 - If $(n_1, n_2, n_3) = (2, 30, 30)$ then $\alpha_1 > \beta_{3,1}$ (pessimistic, conservative)
 - If $(n_1, n_2, n_3) = (2, 3, 4)$, find $\alpha_1 < \beta_{3,1}$ (optimistic, **misleading!**)

CONNECTING SIGNIFICANCE AND FAILURE

Upon terminating a node-level discovery, a significance level will be reported.
How does it relate to the true (conditional) probability of error?

RELATIONSHIP NOT SIMPLE!

- If $k = K$ found, then $\alpha_K > \beta_{K,K} = 0$, so algorithm is conservative
- If $k = K - 1$, also uses $K' = K$, and $\alpha_{K-1} = t_{K,K-1}^{K-1} = \beta_{K,K-1}$, ideal!
- If $k < K - 1$, can go either way, e.g. with $K = 3$ and $k = 1$:
 - If $(n_1, n_2, n_3) = (2, 30, 30)$ then $\alpha_1 > \beta_{3,1}$ (pessimistic, conservative)
 - If $(n_1, n_2, n_3) = (2, 3, 4)$, find $\alpha_1 < \beta_{3,1}$ (optimistic, **misleading!**)

CONNECTING SIGNIFICANCE AND FAILURE

Upon terminating a node-level discovery, a significance level will be reported.
How does it relate to the true (conditional) probability of error?

RELATIONSHIP NOT SIMPLE!

- If $k = K$ found, then $\alpha_K > \beta_{K,K} = 0$, so algorithm is conservative
- If $k = K - 1$, also uses $K' = K$, and $\alpha_{K-1} = t_{K,K-1}^{K-1} = \beta_{K,K-1}$, ideal!
- If $k < K - 1$, can go either way, e.g. with $K = 3$ and $k = 1$:
 - If $(n_1, n_2, n_3) = (2, 30, 30)$ then $\alpha_1 > \beta_{3,1}$ (pessimistic, conservative)
 - If $(n_1, n_2, n_3) = (2, 3, 4)$, find $\alpha_1 < \beta_{3,1}$ (optimistic, **misleading!**)

So can we know anything for sure in practice?

FAILURE BOUNDS

Problem: can we design the $\{n_k\}$ to ensure that $\alpha_k \leq \beta_{K,k}$ for all possible K ?

FAILURE BOUNDS

Problem: can we design the $\{n_k\}$ to ensure that $\alpha_k \leq \beta_{K,k}$ for all possible K ?

Yes !

FAILURE BOUNDS

Problem: can we design the $\{n_k\}$ to ensure that $\alpha_k \leq \beta_{K,k}$ for all possible K ?

Solution:

Step 1: Design the $\{\alpha_k\}$. Easy to show that $\alpha_k \geq t_{K,k}^k \geq t_K^k, \implies$

$$\beta_{K,k} \leq \beta_{K,1} = \sum_{l=1}^{K-1} t_K^l \leq \sum_{l=1}^{K-1} \alpha_l .$$

Note if for example $\alpha_k = \alpha$ for all k , not useful! If set $\alpha_k = \alpha_1 r^{k-1}$:

$$\beta_{K,k} \leq \sum_{l=1}^{K-1} \alpha_l < \frac{\alpha_1}{1-r} .$$

To achieve failure target β^* , select a convenient α_1 , then set $r = 1 - \frac{\alpha_1}{\beta^*}$.

Step 2: From $\{\alpha_k\}$, the corresponding $\{n_k\}$ can be calculated recursively.

Result: $\beta_{K,k} \leq \beta^*$ regardless of the unknown K !

FAILURE BOUNDS

Problem: can we design the $\{n_k\}$ to ensure that $\alpha_k \leq \beta_{K,k}$ for all possible K ?

Solution:

Step 1: Design the $\{\alpha_k\}$. Easy to show that $\alpha_k \geq t_{K,k}^k \geq t_K^k, \implies$

$$\beta_{K,k} \leq \beta_{K,1} = \sum_{l=1}^{K-1} t_K^l \leq \sum_{l=1}^{K-1} \alpha_l .$$

Note if for example $\alpha_k = \alpha$ for all k , not useful! If set $\alpha_k = \alpha_1 r^{k-1}$:

$$\beta_{K,k} \leq \sum_{l=1}^{K-1} \alpha_l < \frac{\alpha_1}{1-r} .$$

To achieve failure target β^* , select a convenient α_1 , then set $r = 1 - \frac{\alpha_1}{\beta^*}$.

Step 2: From $\{\alpha_k\}$, the corresponding $\{n_k\}$ can be calculated recursively.

Result: $\beta_{K,k} \leq \beta^*$ regardless of the unknown K !

EXAMPLE WITH $\beta^* = 0.05$

k	1	2	3	4	5	6	7	8	9
n_k	9	17	24	33	42	51	60	70	81
k	10	11	12	13	14	15	16	...	
n_k	91	102	113	125	136	148	161	...	

TWO KINDS OF STATISTICAL GUARANTEES

Node level guarantees defined by $\{n_k\}$. Can be set in two ways:

- Via failure probability bound: β^*
- Via chosen significance levels: $\alpha_k = \alpha_{\text{node}}$

Graph level guarantees must also bound max size of multipath:

$$\beta_{\text{all}}^* = 1 - \prod_i^{\max} (1 - \beta_{k_i}^*)$$

GROUND TRUTH

TRUE TOPOLOGY UNAVAILABLE

- Network operators don't like sharing
- Network operators don't know (atypical behaviours difficult to pin down)
- Routing is constantly changing..

OUR SURROGATE GROUND TRUTH

- Use MDA with very high β^* to discovery 'everything we can'
- Use same data set to evaluate smaller β^*
 - interfaces found guaranteed to be in ground truth set
 - avoids non-stationarity issue
- Provides well defined assessment of algorithm

GROUND TRUTH

TRUE TOPOLOGY UNAVAILABLE

- Network operators don't like sharing
- Network operators don't know (atypical behaviours difficult to pin down)
- Routing is constantly changing..

OUR SURROGATE GROUND TRUTH

- Use MDA with very high β^* to discovery 'everything we can'
- Use same data set to evaluate smaller β^*
 - interfaces found guaranteed to be in ground truth set
 - avoids non-stationarity issue
- Provides well defined assessment of algorithm

DATA

PROBES

- UDP Probes
- TTL = 36
- Inter-Probe interval: 50ms
- Retransmissions:
 - up to 3 attempts
 - timeout 2 seconds

ROUTES

- Source: UPMC Paris Universit as
 - 44% of routes traverse a LB (prior work gives 23% – 80%)
- Destinations: 5000 randomly selected responsive to pings
- Show 1 of 4 rounds taken over 3 weeks

DATA

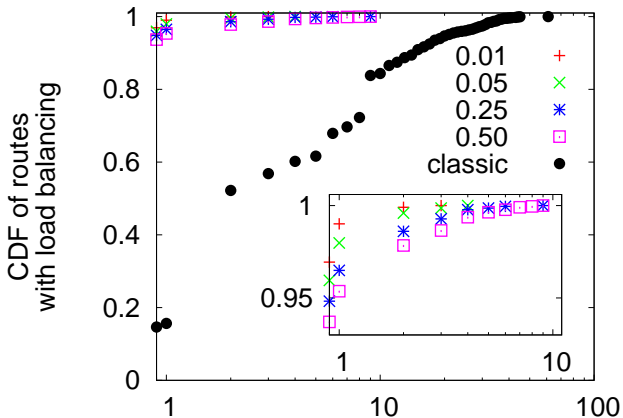
PROBES

- UDP Probes
- TTL = 36
- Inter-Probe interval: 50ms
- Retransmissions:
 - up to 3 attempts
 - timeout 2 seconds

ROUTES

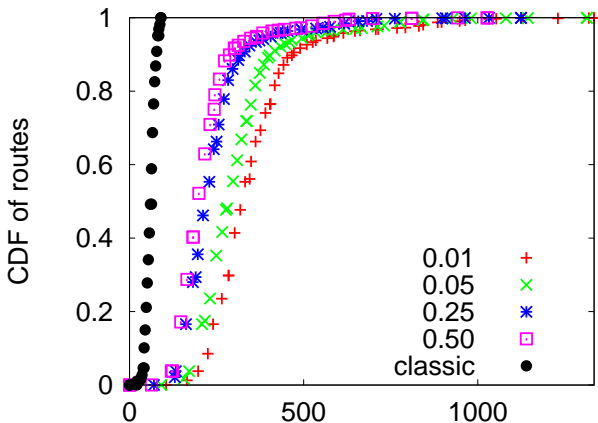
- Source: UPMC Paris Universit as
 - 44% of routes traverse a LB (prior work gives 23% – 80%)
- Destinations: 5000 randomly selected responsive to pings
- Show 1 of 4 rounds taken over 3 weeks

LINK DISCOVERY (LB PATHS ONLY)



- Classic Traceroute inadequate: misses more than 10 links 25% of time!
- Failure bound not tight, $\beta^* = 0.5$ finds full multipath 94% of time

PROBING OVERHEAD (ALL PATHS)



- MDA overhead: flow key control, statistical control, retransmissions
- 20 hop route:
 - Classic Traceroute: mean 61, worst 90
 - MDA ($\beta^* = 0.5$): 216, worst 1027
 - MDA ($\beta^* = 0.01$): 348, worst 1334

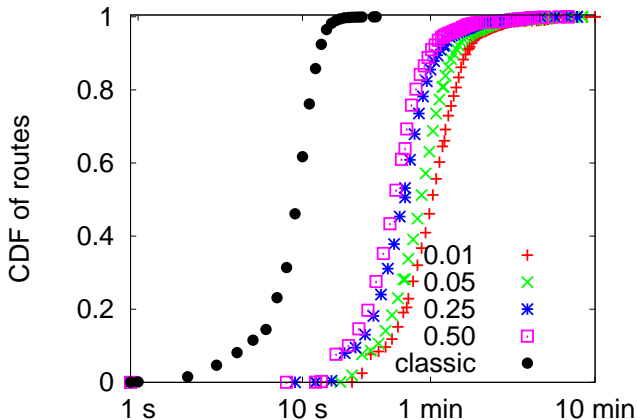
CONCLUSIONS

- Classic Traceroute totally inadequate for multipath discovery
- MDA: an **adaptive** algorithm to find LB generated multipaths
- For the first time: **can measure multipath with confidence** :
 - Failure probabilities β^* at node and graph level
 - universal bound when K unknown
 - with controlled error for target (known) topologies
 - Significance levels α_{node} at node level
- **Future**: reduce probe overhead with tighter bounds and and path control

CONCLUSIONS

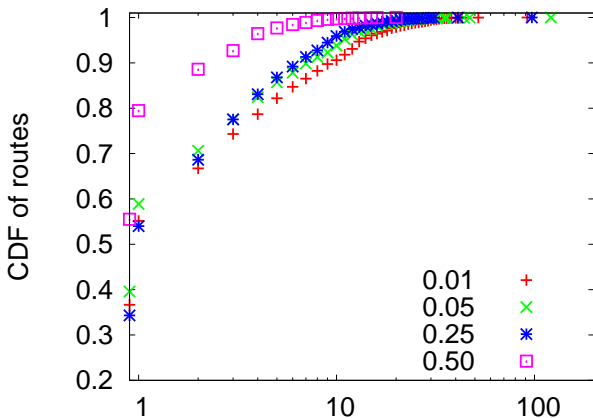
- Classic Traceroute totally inadequate for multipath discovery
- MDA: an **adaptive** algorithm to find LB generated multipaths
- For the first time: **can measure multipath with confidence** :
 - Failure probabilities β^* at node and graph level
 - universal bound when K unknown
 - with controlled error for target (known) topologies
 - Significance levels α_{node} at node level
- **Future**: reduce probe overhead with tighter bounds and and path control

TIME OVERHEAD (ALL PATHS)



- MDA clearly much greater in general
- Parallelism difficult but worth investigating

RETRANSMISSIONS (ALL PATHS, MDA ONLY)



- Retransmissions very common! (45% of routes for $\beta^* = 0.5$)
- Scarcely worth it, most routers reply always, or never..