

Modeling User Submission Strategies on Production Grids

Diane Lingrand, Johan Montagnat, Tristan Glatard

Sophia-Antipolis, FRANCE

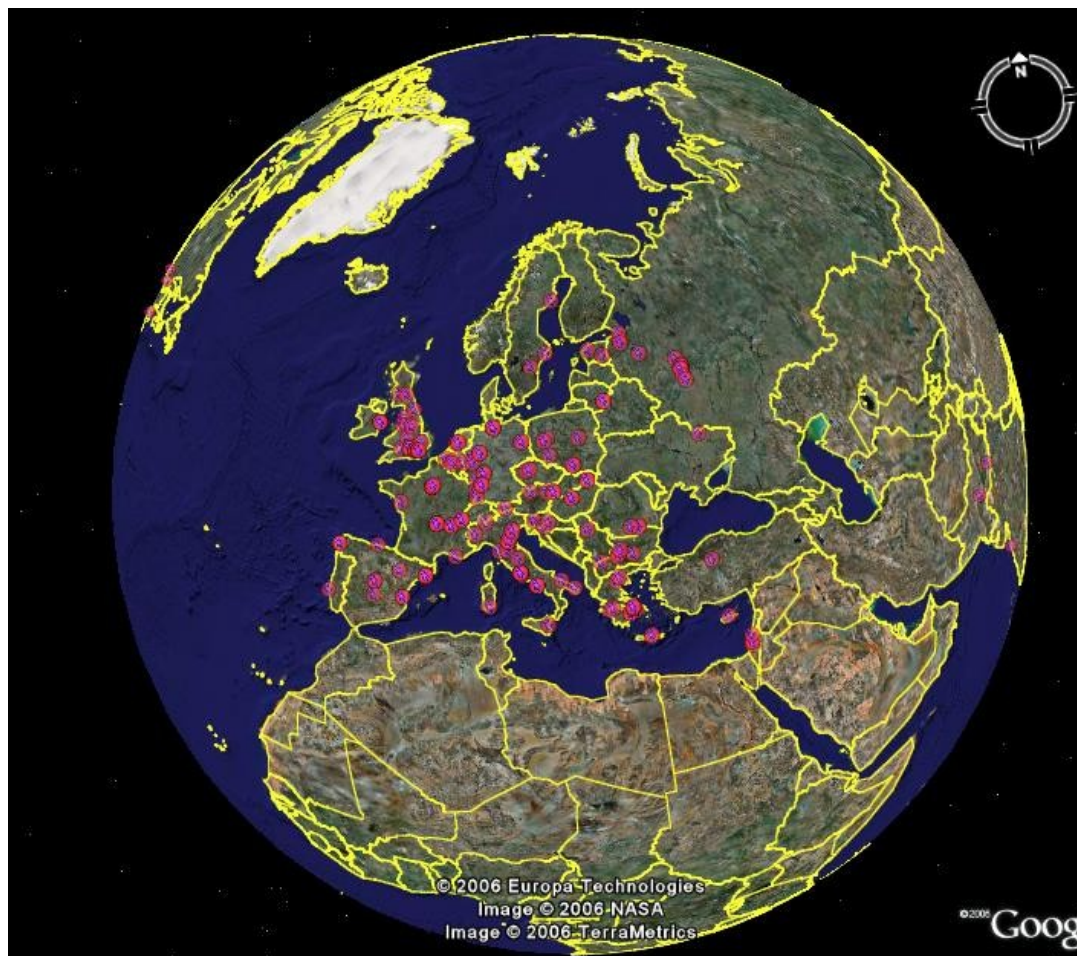
University of Lyon, FRANCE



HPDC 2009
München, Germany

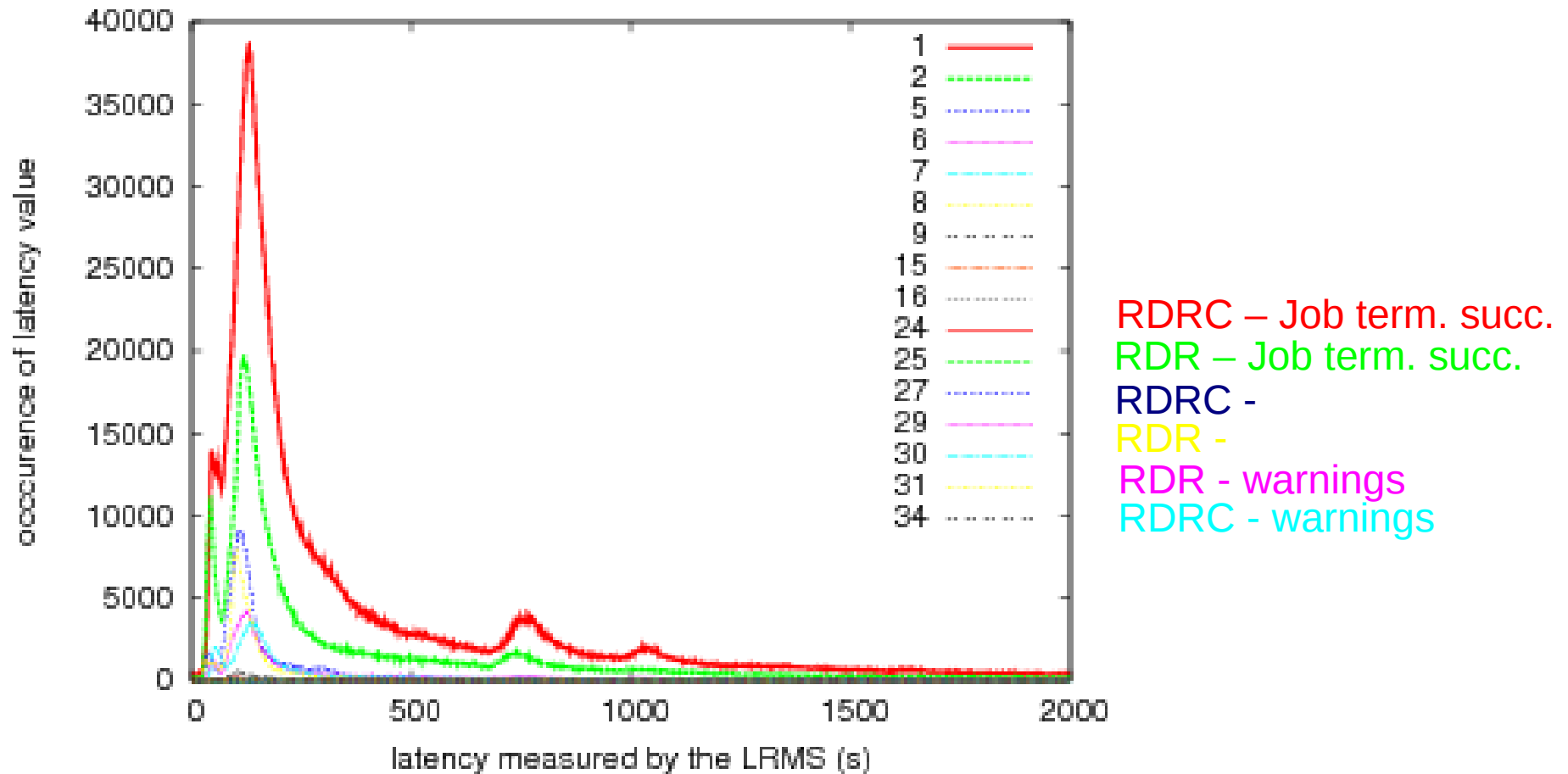
The EGEE production grid

<http://www.eu-egee.org>



- Huge computing power and data storage facility:
 - > 80,000 CPUs
 - > 250 computing centers world-wide
 - > 200,000 jobs/day
 - > 9,000 registered users
- Toll: latency and faults

Variable latencies

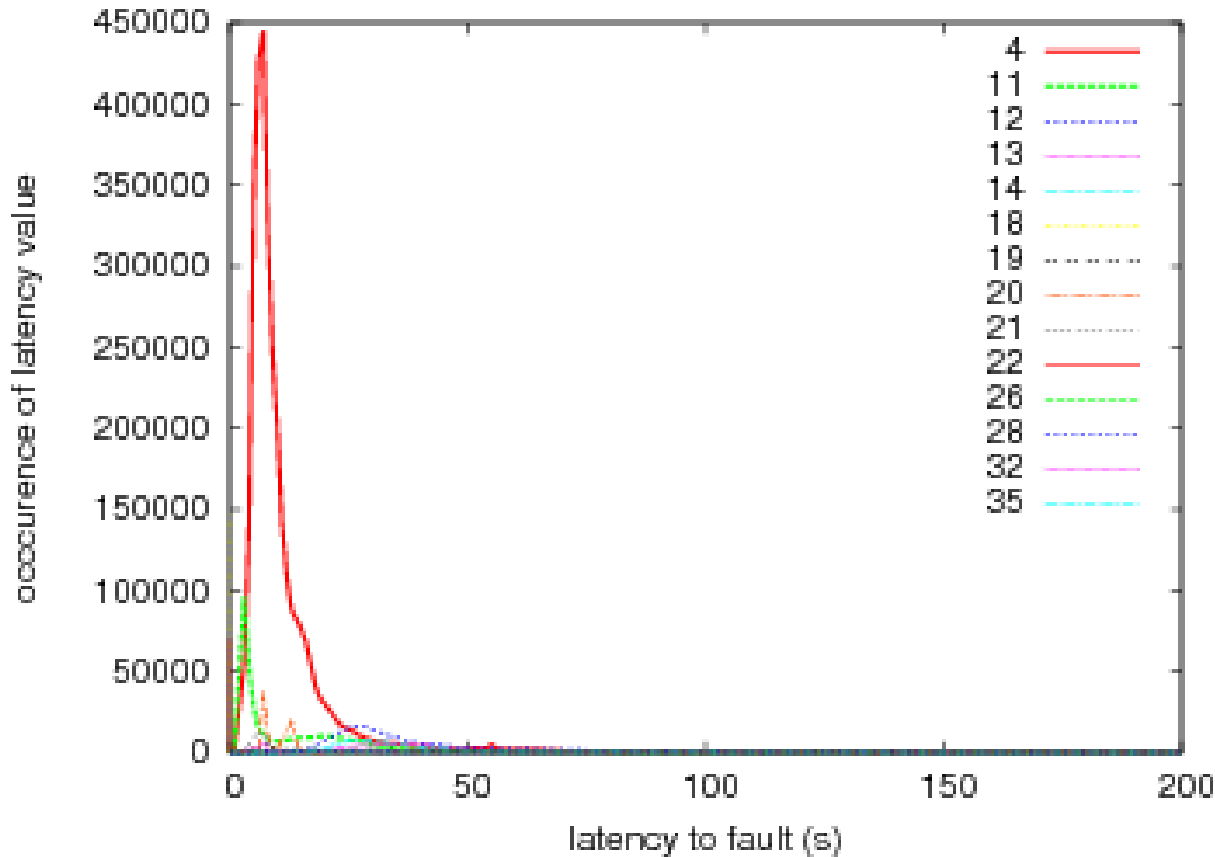


- Heavy-tailed, multimodal

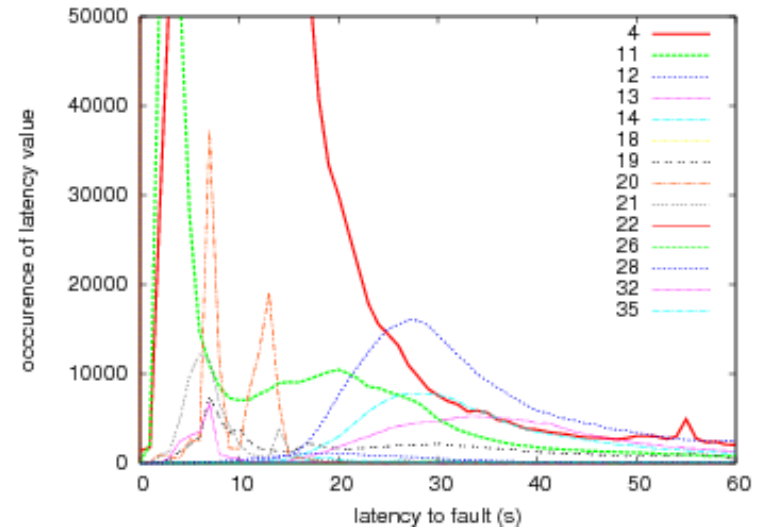
Based on 33 millions of EGEE jobs, 2005-2007

[Lingrand *et al*, JSSPP'09]

Faults



- RA - No compatible resource
- RA -
- RA - Job proxy expired
- RA - cannot retrieve previous matches

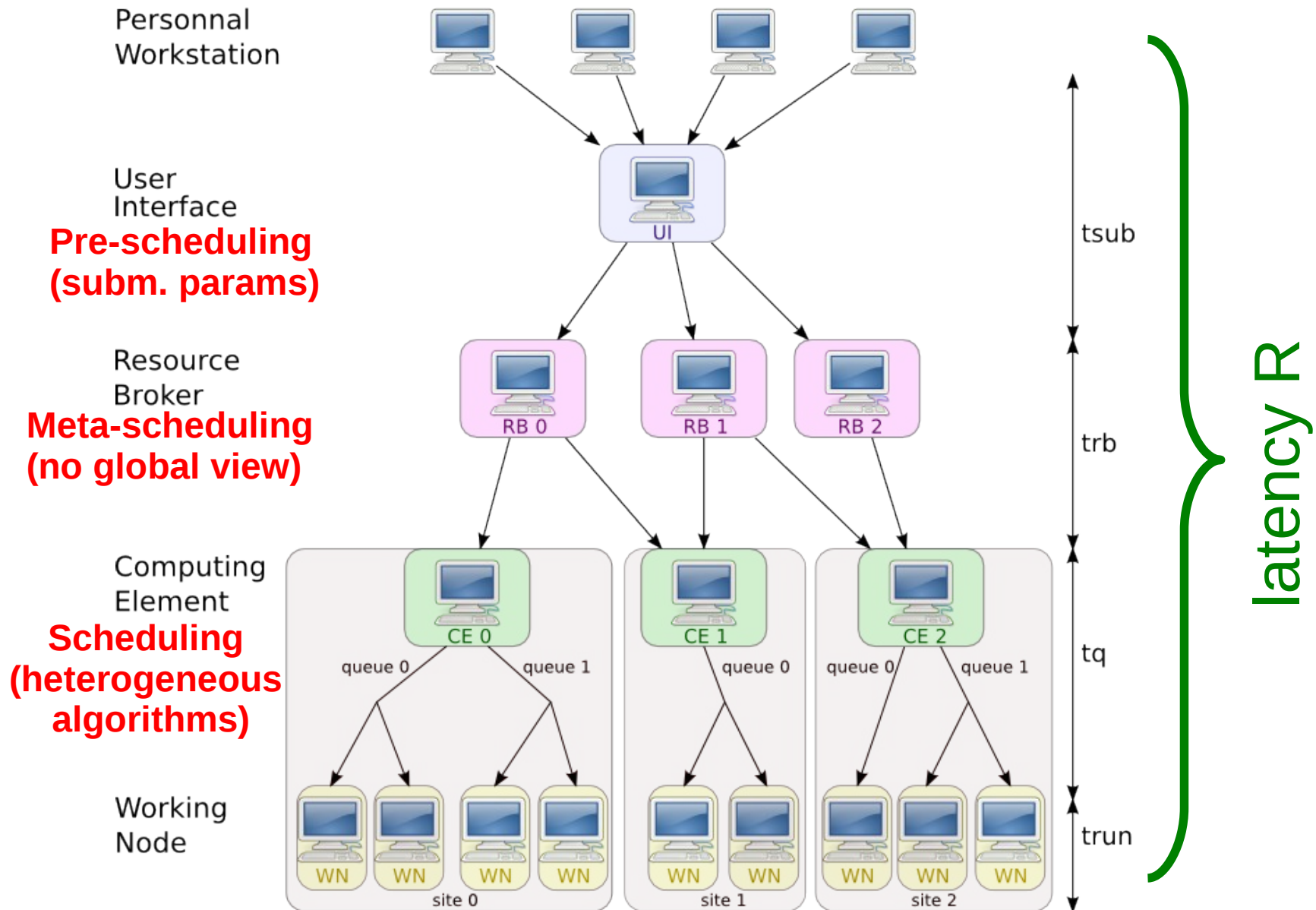


- 35% of faults, various distributions

Based on 33 millions of EGEE jobs, 2005-2007

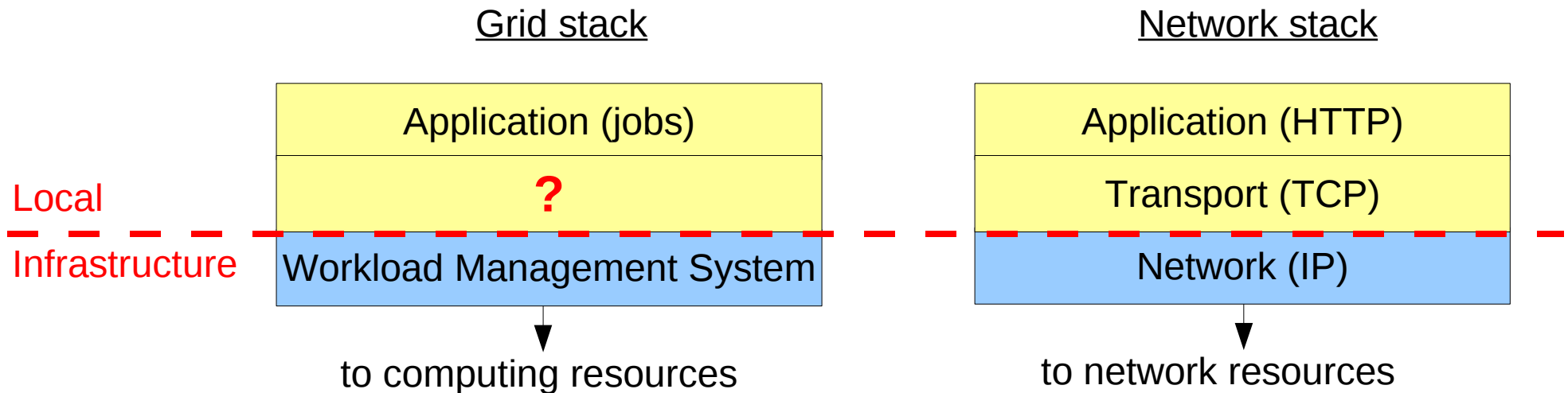
[Lingrand *et al*, JSSPP'09]

EGEE: a complex system



Quality of Service for grid applications

- Assumption: infrastructure only provides best-effort
 - Analogy with TCP/IP model



[Meng *et al* IPDPS'09]

- => Need for user-level submission strategies

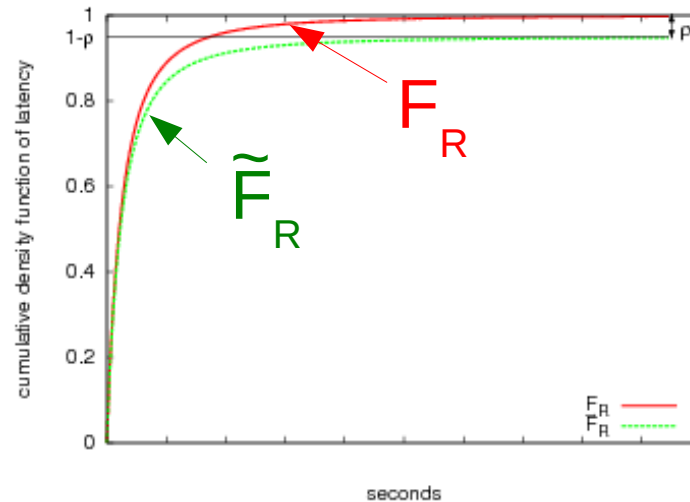
Outline

- Introduction
- Model and evaluation of 3 strategies
 - Single resubmission after timeout [CCGrid'07]
 - Multiple submissions [Casanova, JGC 07]
 - Delayed resubmission
- Cost criterion
- Conclusion

Modeling & evaluation method

- Infrastructure behavior

- Latency R is a random variable with c.d.f F_R
- Some jobs (outliers, fraction ρ) have infinite latency



- Submission strategies modeling

- Expectation & stdev of total latency J w.r.t. R and ρ

- Evaluation

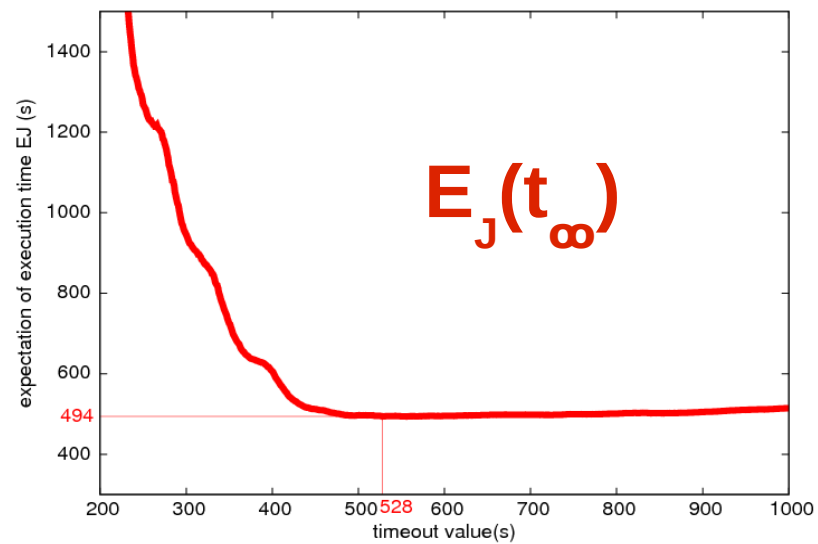
- Based on 11,000 probe jobs (Sept. 06 – Feb 08)

Single resubmission: modeling

- Job timed-out => canceled and resubmitted
- Expectation of total latency time:

$$E_J(t_\infty) = \frac{1}{\tilde{F}_R(t_\infty)} \int_0^{t_\infty} (1 - \tilde{F}_R(u)) du$$

- E_J has a min $\Leftrightarrow \tilde{F}_R$ is heavy-tailed



Single resubmission: evaluation

- Total latencies

week number	Without outliers	With outliers	E_J	Without outliers		
	mean < 10 ⁴ s	mean with 10 ⁴ s		σ_R < 10 ⁴ s	σ_J	$\Delta\sigma$
2006-IX	570s	1042s	471s	886s	331s	-63%
2007/08	469s	2089s	500s	723s	358s	-51%
2007-36	446s	2739s	510s	748s	370s	-51%
2007-37	506s	3639s	617s	848s	486s	-43%
2007-38	447s	2739s	531s	682s	399s	-42%
2007-39	489s	3533s	596s	741s	482s	-35%
2007-50	660s	2341s	628s	1046s	475s	-55%
2007-51	478s	1716s	517s	510s	353s	-31%
2007-52	443s	1685s	476s	582s	334s	-43%
2007-53	449s	1977s	482s	678s	330s	-51%
2008-01	434s	1678s	499s	317s	339s	+07%
2008-02	418s	1568s	441s	547s	278s	-49%
2008-03	538s	1484s	419s	1196s	269s	-78%

- Conclusions

- Manages to filter out outliers
- Reduces standard-deviation

Multiple resubmission: modeling

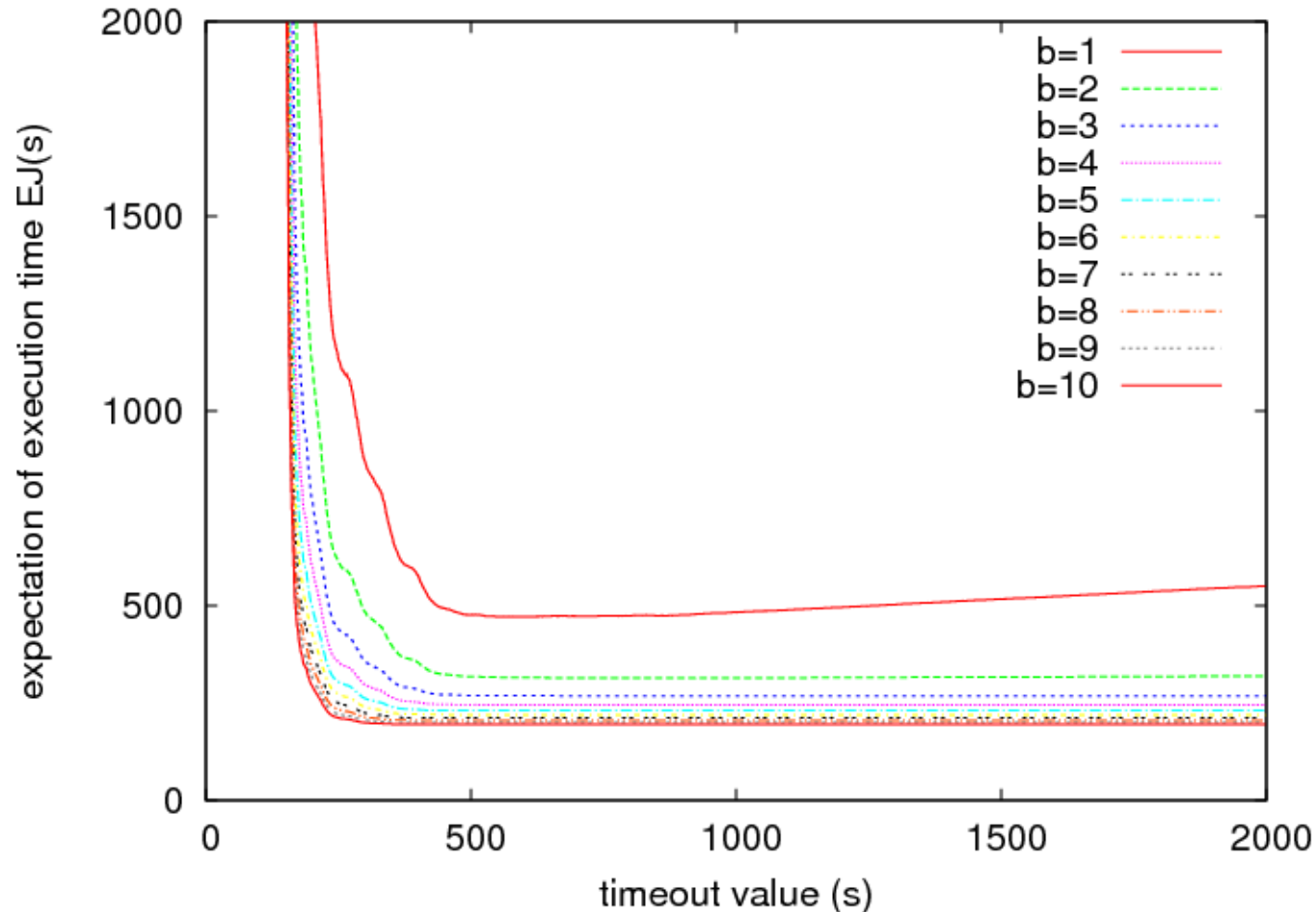
- Submit b copies of the job
 - With timeout on the collection
- Expectation of total latency time

$$E_J(t_\infty) = \frac{1}{\tilde{F}_R(t_\infty)} \int_0^{t_\infty} (1 - \tilde{F}_R(u)) du$$

$\tilde{F}_R(t)$ replaced by $1 - \underbrace{(1 - \tilde{F}_R(t))^b}_{\substack{\text{b jobs have latency } > t \\ \text{At least 1 job has latency } < t}}$

$$E_J(t_\infty) = \frac{1}{1 - (1 - \tilde{F}_R(t_\infty))^b} \int_0^{t_\infty} (1 - \tilde{F}_R(u))^b du$$

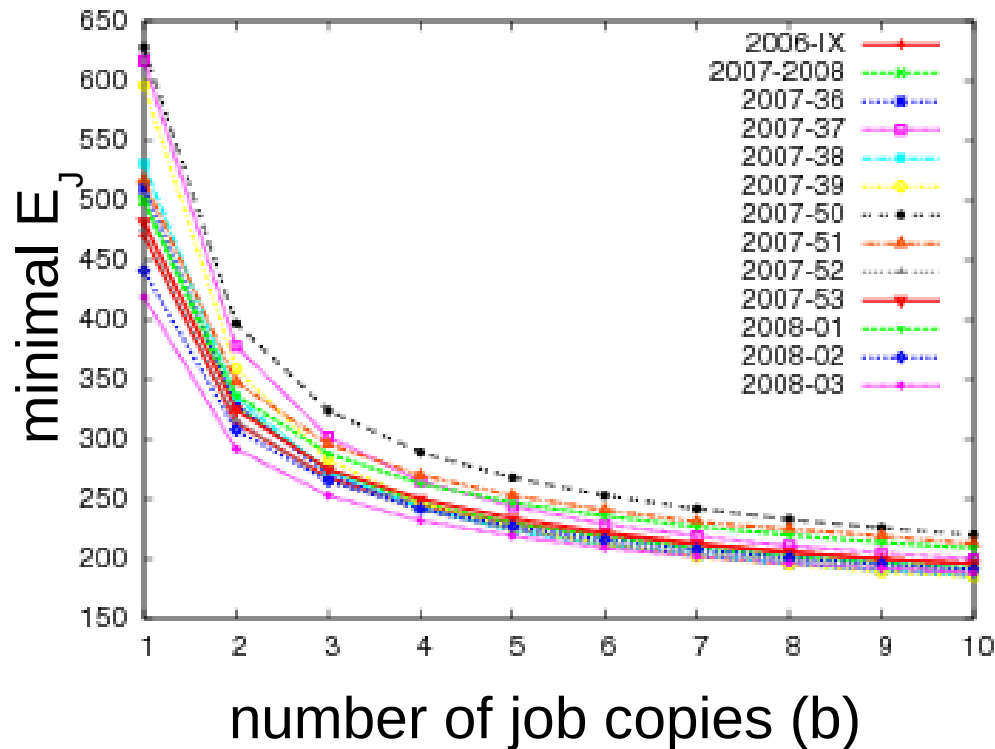
Multiple resubmission: evaluation



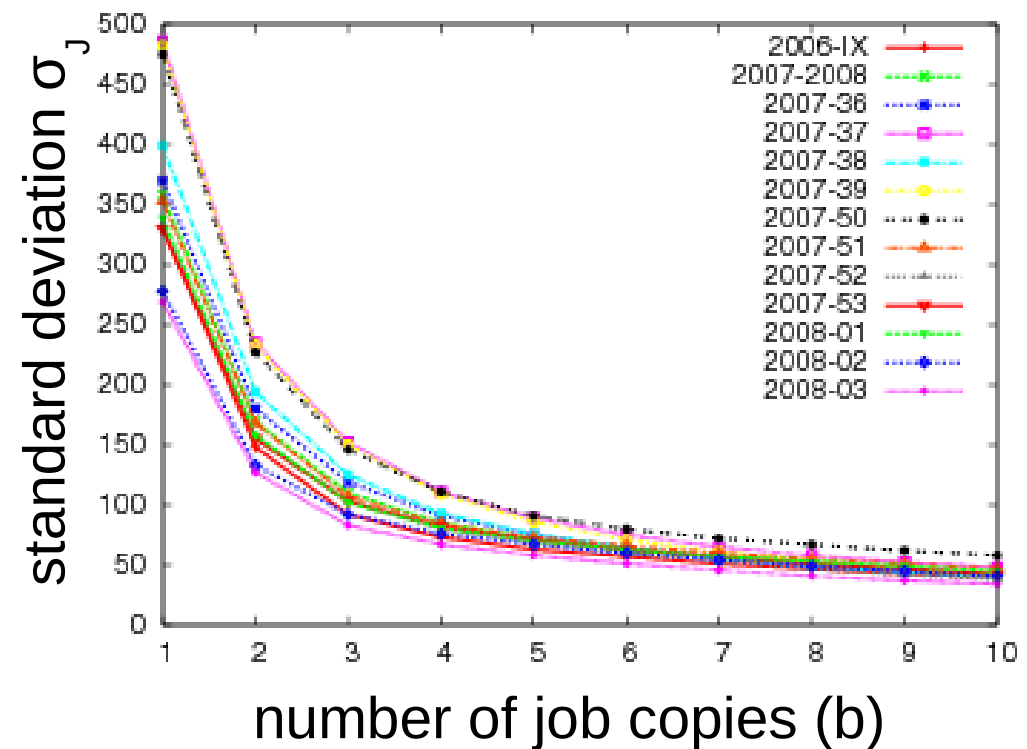
- E_J has a minimum for all values of b
- Slope after minimum decreases as b increases

Multiple resubmission: evaluation

Mean



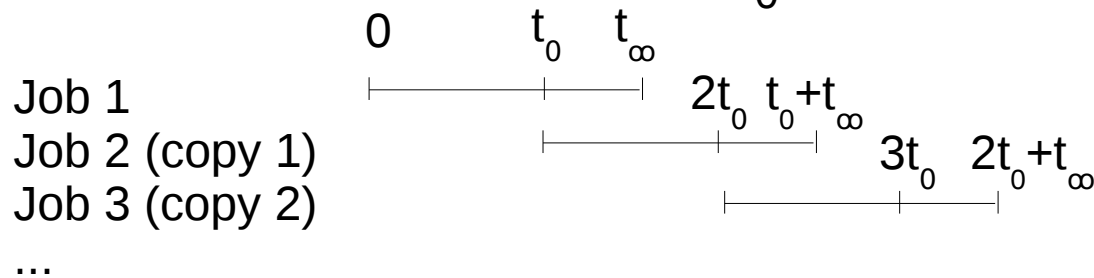
Stdev



- Strong improvements of mean and stdev
- Goes smoother as b increases

Delayed resubmission strategy

- Goal: limit the number of simultaneous job copies
 - Job submitted at time t with timeout t_∞
 - Copy submitted at $t+t_0$



Constraints:

$$0 < t_0 < t_\infty$$

$$t_\infty < 2t_0$$

- Expectation of total latency:

$$\begin{aligned}
 E_J(t_0, t_\infty) = & \frac{1}{\tilde{F}_R(t_\infty)} \int_0^{t_\infty} u \tilde{f}_R(u) du + \frac{\tilde{F}_R(t_0)}{\tilde{F}_R(t_\infty)} \int_0^{t_\infty - t_0} u \tilde{f}_R(u) du \\
 & + \frac{t_0}{\tilde{F}_R(t_\infty)} + t_0 \frac{\tilde{F}_R(t_\infty - t_0)}{\tilde{F}_R(t_\infty)} + t_0 \frac{\tilde{F}_R(t_0) \tilde{F}_R(t_\infty - t_0)}{\tilde{F}_R^2(t_\infty)} - t_0 + \int_0^{t_\infty - t_0} u \tilde{f}_R(u) du \\
 & - \frac{t_0}{\tilde{F}_R(t_\infty)^2} \int_0^{t_\infty - t_0} \tilde{f}_R(u + t_0) \cdot \tilde{f}_R(u) du - \frac{1}{\tilde{F}_R(t_\infty)} \int_0^{t_\infty - t_0} u \tilde{f}_R(u + t_0) \cdot \tilde{f}_R(u) du
 \end{aligned}$$

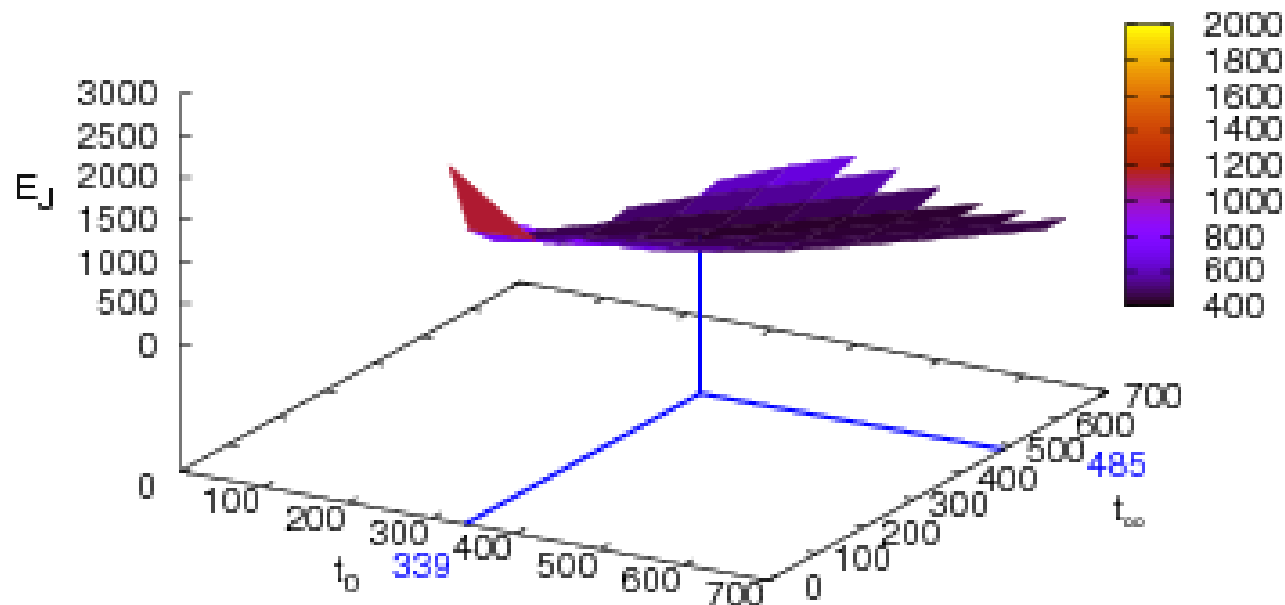
Delayed resubmission strategy

- Minimal value of E_J

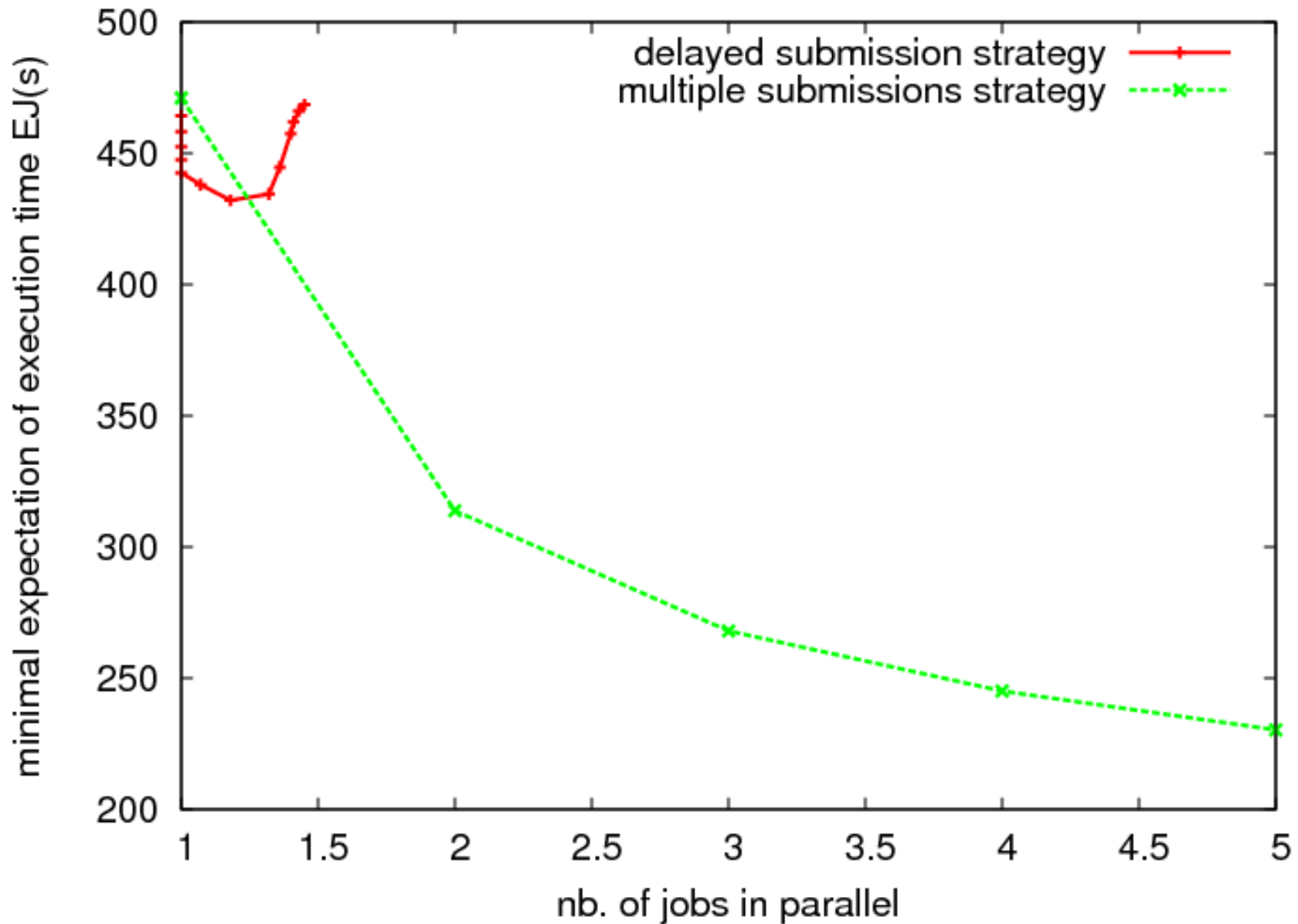
Single resubmission: $E_J = 471s$

Multiple resubmissions ($b=2$): $E_J = 314s$

Delayed resubmission: $E_J = 431s$



Total latency of delayed VS mult. subm.

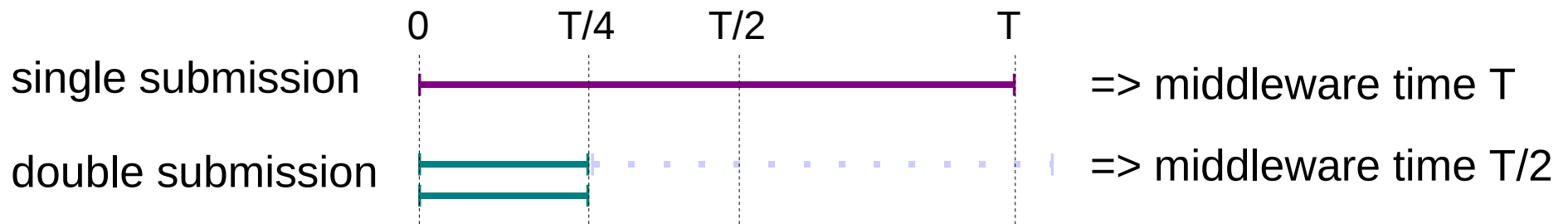


Outline

- Introduction
- Model and evaluation of 3 strategies
 - Single resubmission after timeout
 - Multiple submissions
 - Delayed resubmission
- Cost criterion
- Conclusion

Cost of the strategies

- Based on total middleware time



- When $N_{//}$ copies are submitted in parallel:

$$\Delta_{cost} = N_{//} * \frac{E_J(\text{with } N_{//})}{E_J(\text{with } b = 1)}$$

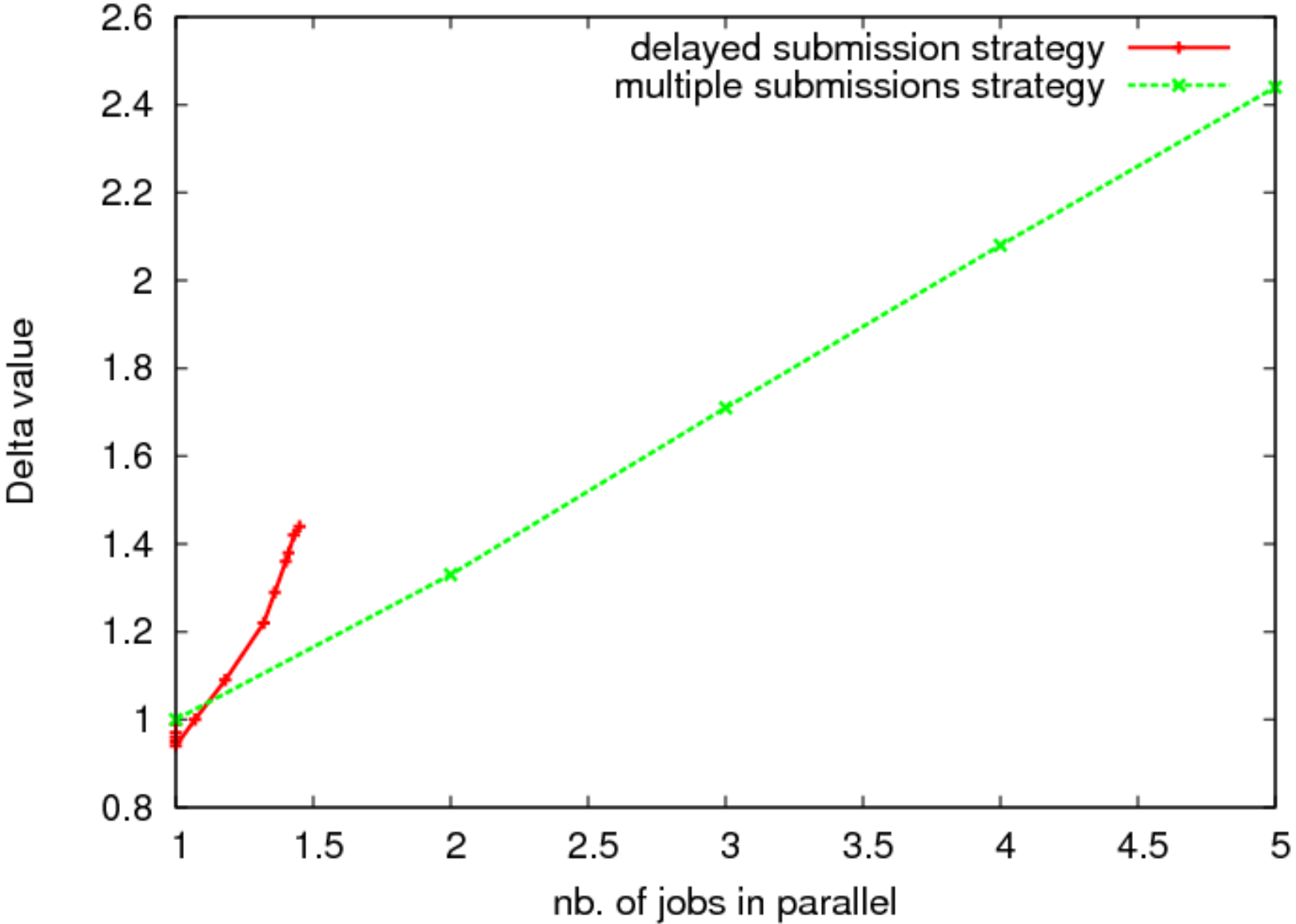
Cost of delayed VS multiple submission

$N_{//}$	$\frac{t_{\infty}}{t_0}$	E_J	Δ_{cost}	$N_{//}$	E_J	Δ_{cost}	$N_{//}$	E_J	Δ_{cost}
1		471s	1	2	314s	1.3	16	180s	6.1
1	1.1	458s	0.97	3	268s	1.7	17	178s	6.4
1	1.15	453s	0.96	4	245s	2.1	18	177s	6.7
1	1.2	447s	0.95	5	230s	2.4	19	175s	7.1
1	<u>1.25</u>	443s	<u>0.94</u>	6	220s	2.8	20	174s	7.4
1.07	1.3	438s	1.00	7	212s	3.1	30	166s	10
1.18	1.4	432s	1.09	8	205s	3.5	40	161s	14
1.32	1.5	434s	1.22	9	200s	3.8	50	158s	17
1.36	1.6	445s	1.29	10	196s	4.2	60	156s	20
1.40	1.7	458s	1.36	11	192s	4.5	70	155s	23
1.41	1.8	462s	1.38	12	189s	4.8	80	154s	26
1.43	1.9	466s	1.42	13	186s	5.1	90	153s	29
1.45	2.0	469s	1.44	14	184s	5.5	100	152s	32
				15	182s	5.8			

delayed resubmission strategy

multiple resubmission strategy

Cost of delayed VS multiple submission



Conclusion

- Need for local user-level submission strategies
 - Modeling
 - Evaluation
- Studied 3 submission strategies
 - Single, multiple, delayed
- Cost metric based on total middleware time
- Conclusions
 - Multiple resubmission reduces total latency at a high cost ($b=2 \Rightarrow \Delta_{\text{cost}} = 1.3$)
 - Delayed strategy improves total latency at a lower cost than single resubmission
- Future work: do it live on real applications

Thank you !

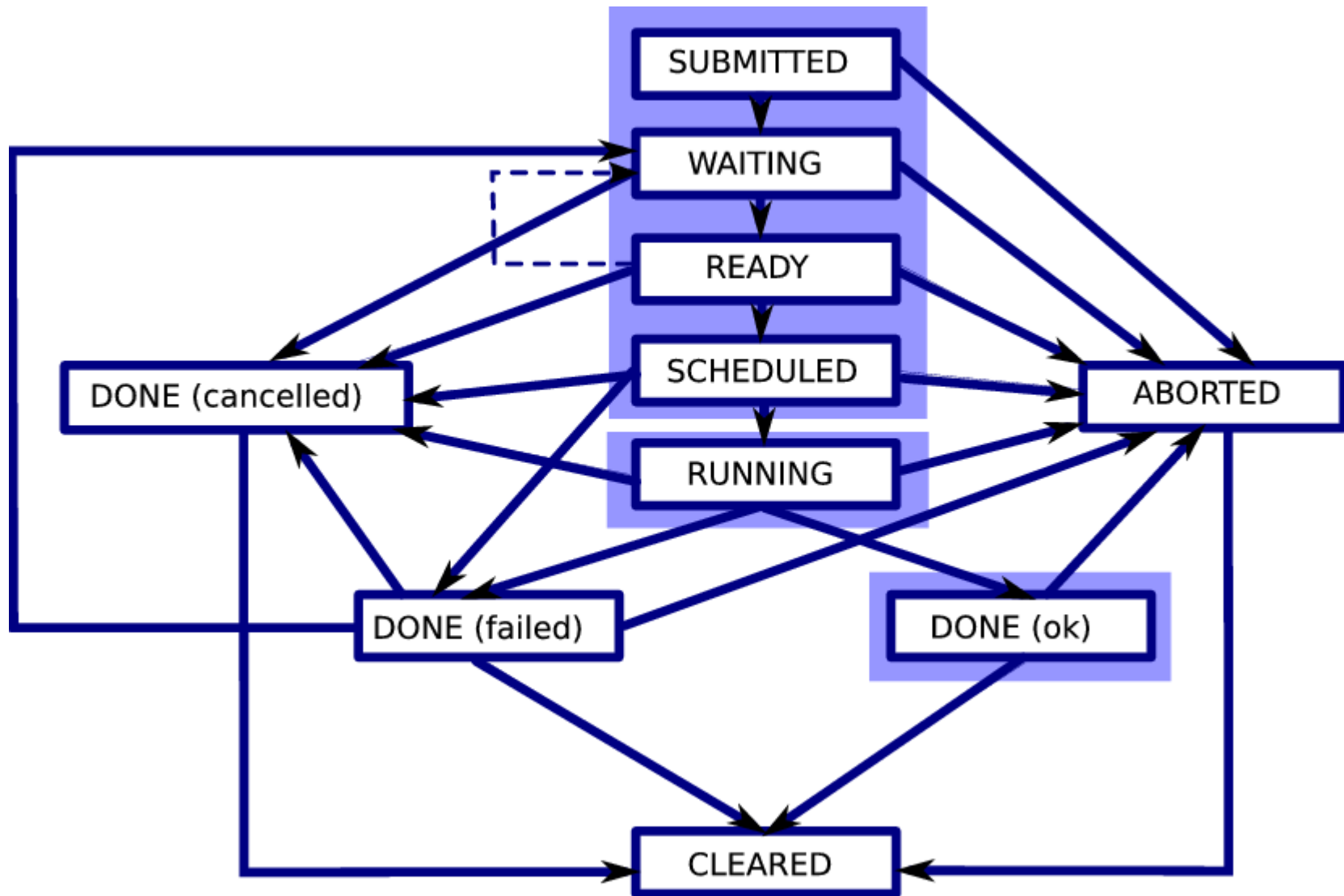
Questions ?

Back slides

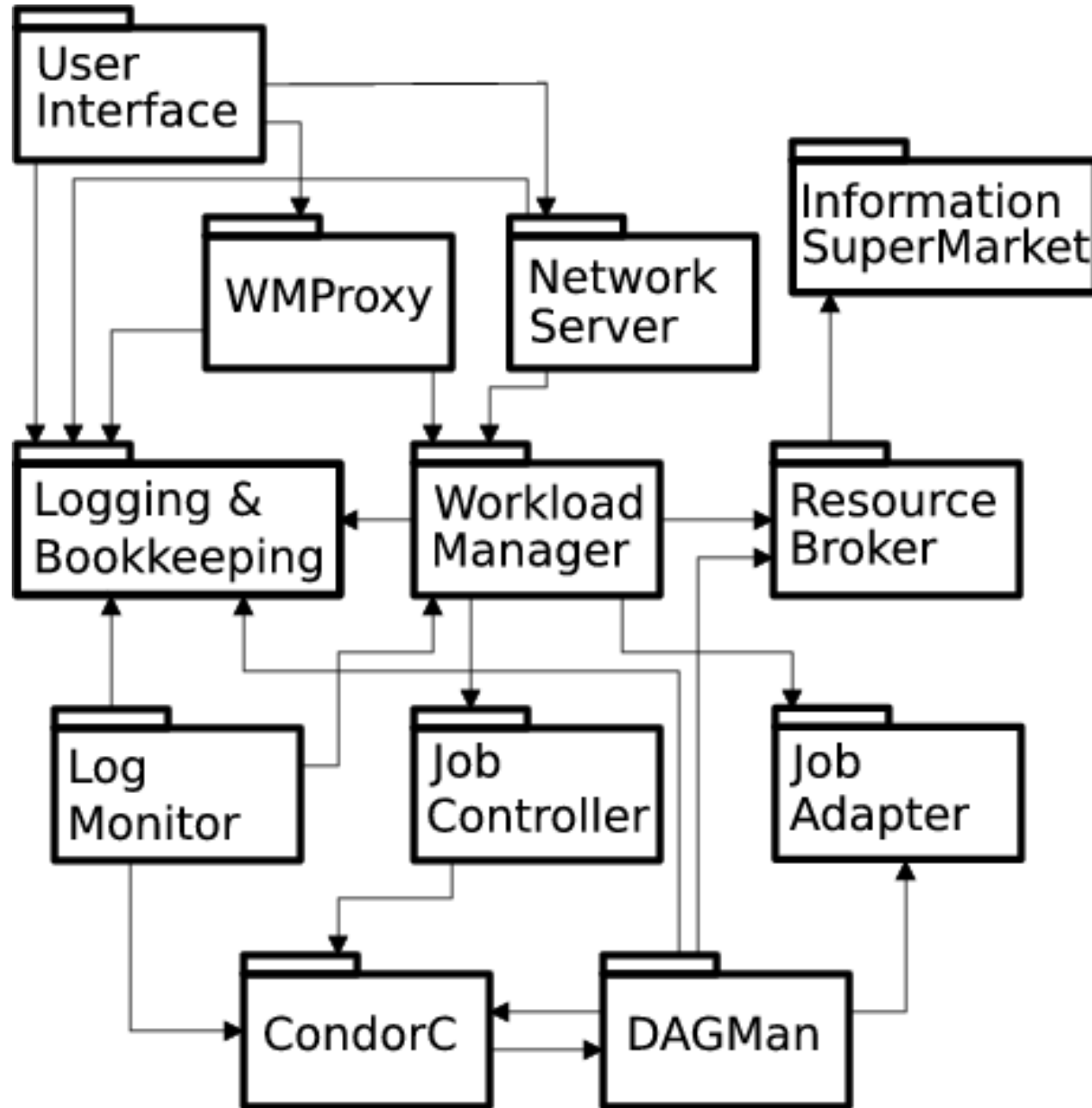
Future work

- Implementations in real applications
 - Need for latency estimates
- Using data from the Grid Observatory
 - <http://www.grid-observatory.org>
 - non sparse data
 - non limited to the biomed VO

Job's life cycle on EGEE



EGEE Workload Management System



Data

- probe jobs: /bin/hostname
- periodically submitted to maintain constant load
- 11,000 traces acquired
 - on the biomed VO
 - between September 2006 and February 2008