



Modelling job allocation where service duration is unknown

Nigel Thomas

University of Newcastle upon Tyne



Contents

- Motivation
- TAG
- Aims
- The basic model
- Some initial results
- Phase type distributions
- More results
- Conclusions



Motivation: Grid computing

- The promise of grid computing is seamless integration of multiple resources from multiple domains.
 - Customers are promised the ability to automatically select the best resources for their task.
 - Providers are promised the ability to schedule tasks efficiently according to whatever criteria they specify.
- Currently the reality is somewhat different!
- This is due (in part) to the difficulty of finding the necessary information.
- Therefore an important research topic is determining the best strategies to employ in grid scheduling and brokerage with limited available information.



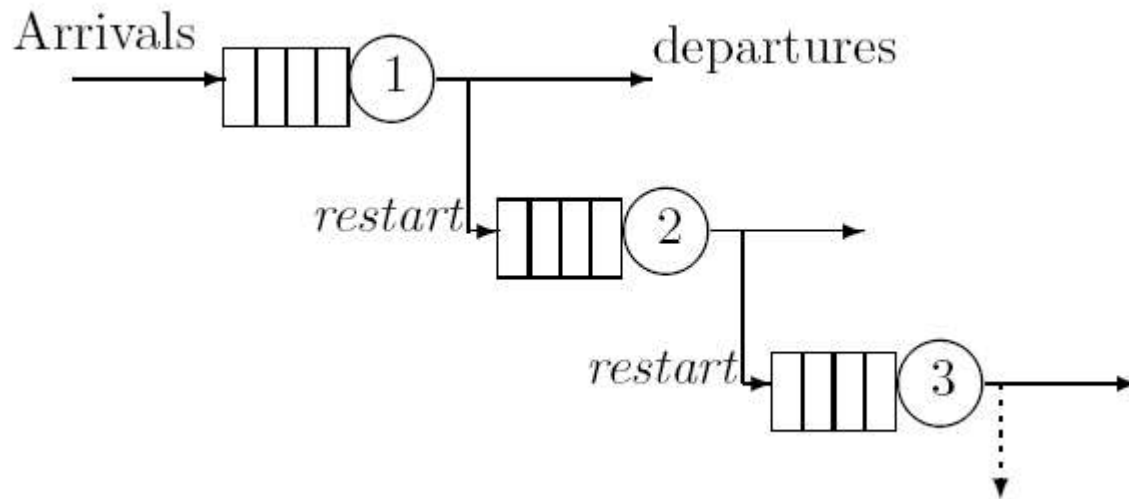
Motivation: fair use of resources

- Suppose there are N jobs waiting to be served, but the head job is much larger than all the others, then all the jobs must wait a long time.
- If we have a choice of servers/queues then we would like to allocate jobs queues according to their size to prevent (or reduce) this problem.
- However, estimating job duration is notoriously costly and unreliable, so what can we do?
 - Pull jobs from a central resource.
 - Assign jobs to the service centre with the shortest queue.
 - Assign jobs to service centres on a round robin basis.
 - Assign jobs to service centres randomly according to probabilities based on service capacity.



TAG (Mor Harchol-Balter, 2002)

- Seeks to reduce variance at each queue.
- Allows shorter jobs to proceed quickly, but delays longer jobs more.
- Employs a progressive filtering mechanism.
- Stops jobs that execute for too long and re-queues them at the subsequent node.





Aims of this paper

- To model the TAG scheme in PEPA.
- To investigate different service distributions (initially negative exponential and Cox).
- To investigate finite capacity queues, and hence consider job loss.



The model

- A 2 node system modelled in PEPA.
- The deterministic timeout is modelled with an Erlang distribution.
- Queues are modelled along conventional state based lines (queue-centric).
- The second queue employs an Erlang repeat service followed by the exponential remainder.

$$Node_1 \begin{array}{c} \boxtimes \\ \{timeout\} \end{array} Node_2$$



Node 1

$$Queue1_0 \stackrel{def}{=} (arrival, a).Queue1_1$$

$$Queue1_i \stackrel{def}{=} (arrival, a).Queue1_{i+1} + (service1, \top).Queue1_{i-1} \\ + (timeout, \top).Queue1_{i-1} + (tick1, \top).Queue1_i, 1 \leq i < K_1$$

$$Queue1_{K_1} \stackrel{def}{=} (timeout, \top).Queue1_{K_1-1} + (service1, s).Queue1_{K_1-1} \\ + (tick1, \top).Queue1_{K_1}$$

$$Server1 \stackrel{def}{=} (service1, s).Server1$$

$$Time1_0 \stackrel{def}{=} (timeout, t).Time1_n + (service1, \top).Time1_n$$

$$Time1_i \stackrel{def}{=} (tick1, t).Time1_{i-1} + (service1, \top).Time1_n, 1 \leq i \leq n$$

$$Node_1 \stackrel{def}{=} Server1 \underset{\{service1\}}{\boxtimes} Queue1_0 \left\{ \underset{service1, tick1}{\boxtimes}^{timeout} \right\} Time1_n$$

Node 2



$$Queue2_0 \stackrel{def}{=} (timeout, \top).Queue2_1$$

$$Queue2_i \stackrel{def}{=} (timeout, \top).Queue2_{i+1} + (repeat\ service, \top).Q2_i \\ + (tick2, \top).Queue2_i, 1 \leq i < K_2$$

$$Q2_i \stackrel{def}{=} (timeout, \top).Q2_{i+1} + (service2, s).Queue2_{i-1} \\ + (tick2, \top).Q2_i, 1 \leq i < K_2$$

$$Queue2_{K_2} \stackrel{def}{=} (timeout, \top).Queue2_{K_2} + (repeat\ service, \top).Q2_{K_2} \\ + (tick2, \top).Queue2_{K_2}$$

$$Q2_{K_2} \stackrel{def}{=} (timeout, \top).Q2_{K_2} + (service2, s).Queue2_{K_2-1} \\ + (tick2, \top).Q2_{K_2}$$

$$Server2 \stackrel{def}{=} (service2, s).Server2$$

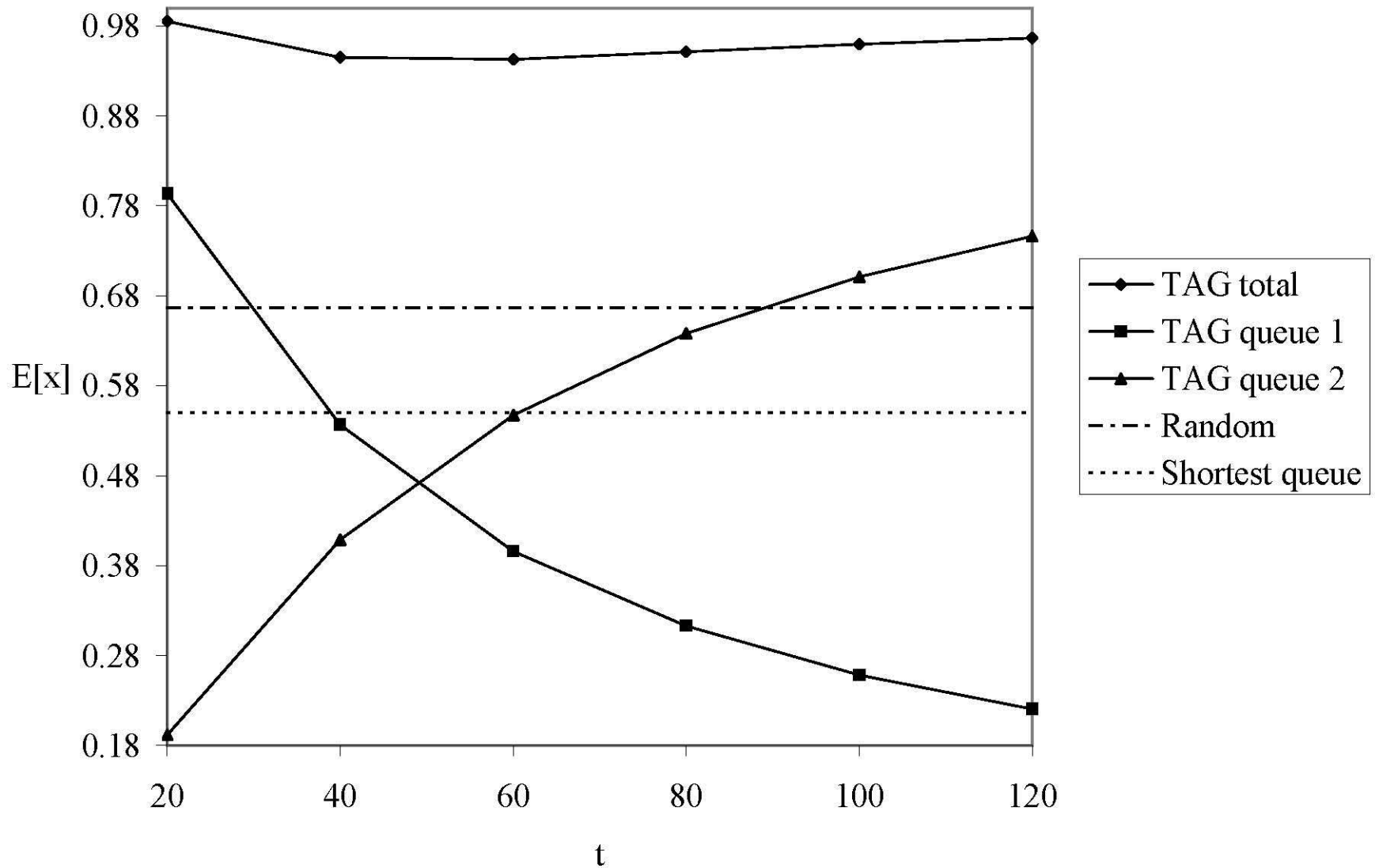
$$Time2_0 \stackrel{def}{=} (repeat\ service, t).(service2, \top).Time2_n$$

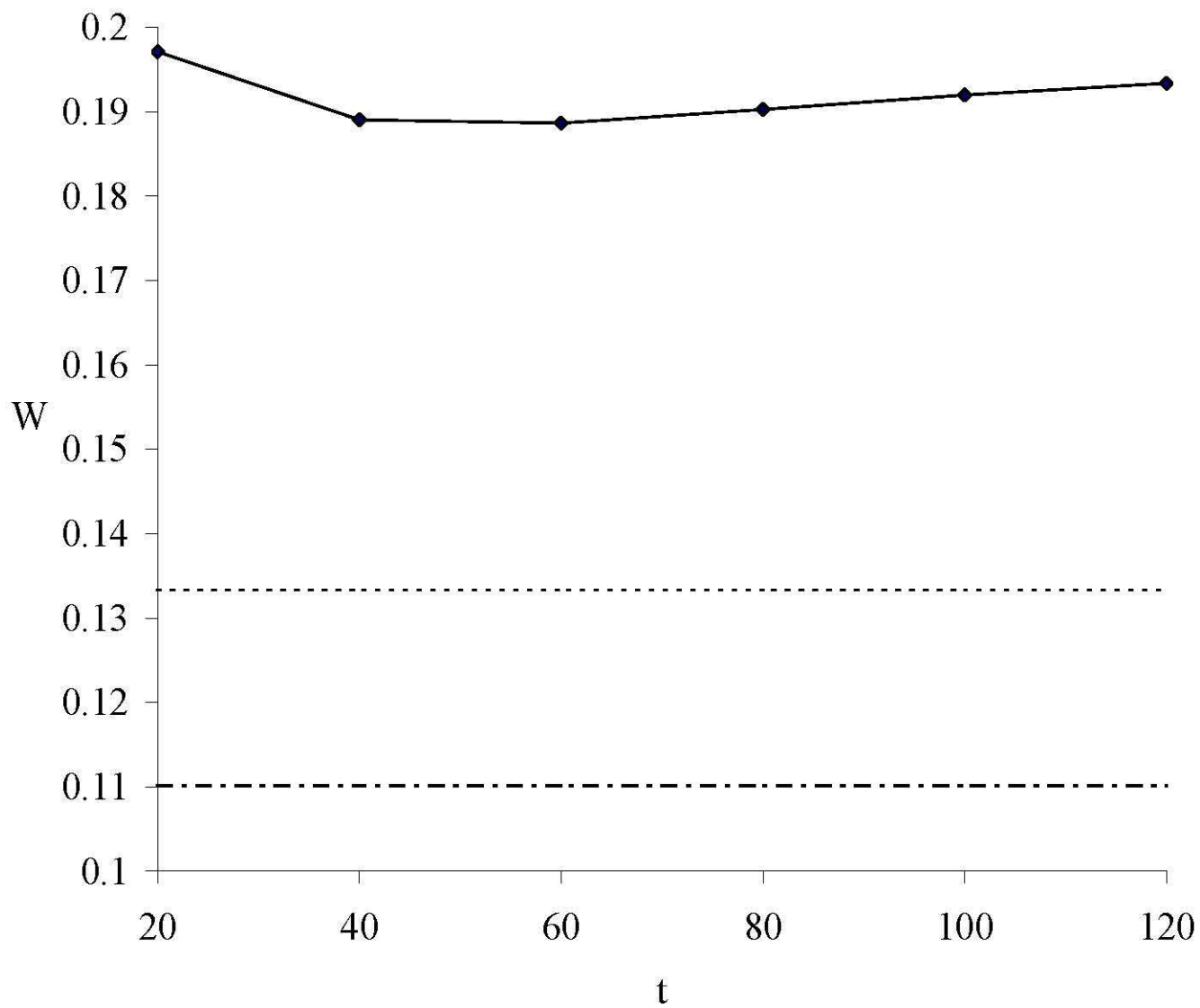
$$Time2_i \stackrel{def}{=} (tick2, t).Time2_{i-1}, 1 \leq i \leq n$$

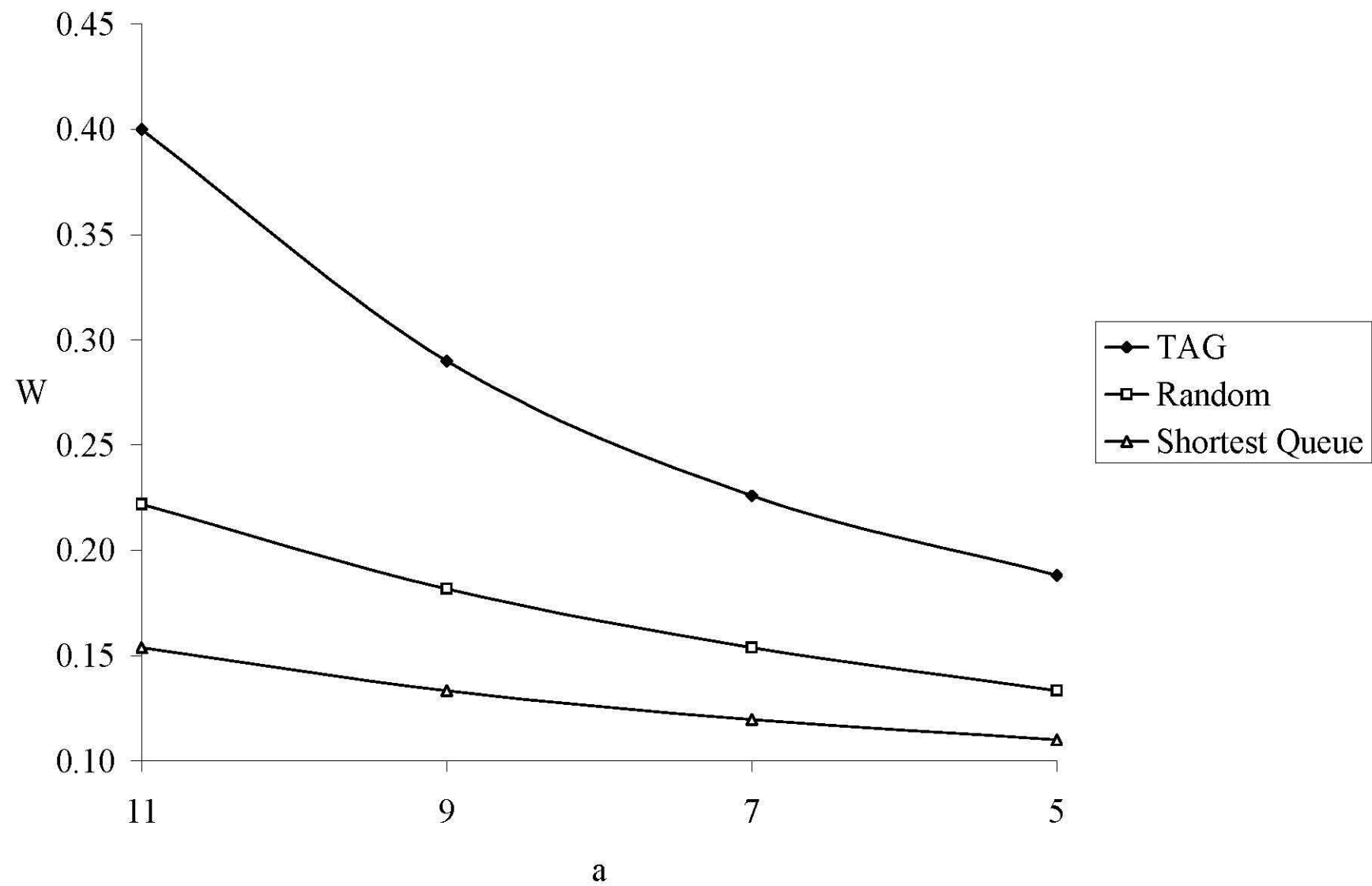
$$Node_2 \stackrel{def}{=} Server2 \boxtimes_{\{service2\}} Queue2_0 \boxtimes_{\{repeat\ service, service2, tick2\}} Time2_n$$



Initial results (1/3)









Phase type distributions

- Consist of multiple exponentials, e.g.
 - Erlang
 - Hyper-exponential

$$\text{Start}H_2 \stackrel{\text{def}}{=} (\text{branch}, \alpha \times \text{fast}).(\text{phase1}, \lambda_1).\text{End} \\ + (\text{branch}, (1 - \alpha) \times \text{fast}).(\text{phase2}, \lambda_2).\text{End}$$

- Cox

$$\text{Start}C_2 \stackrel{\text{def}}{=} (\text{phase1}, \alpha\lambda_1).\text{End} + (\text{phase1}, (1 - \alpha)\lambda_1).(\text{phase2}, \lambda_2).\text{End}$$



Adding Cox distributed service

$$Server1 \stackrel{def}{=} (service1, \alpha s_1).Server1 + (phase1, (1 - \alpha)s_1).Server1' + (timeout, \top).Server1$$

$$Server1' \stackrel{def}{=} (service1, s_2).Server1 + (timeout, \top).Server1$$

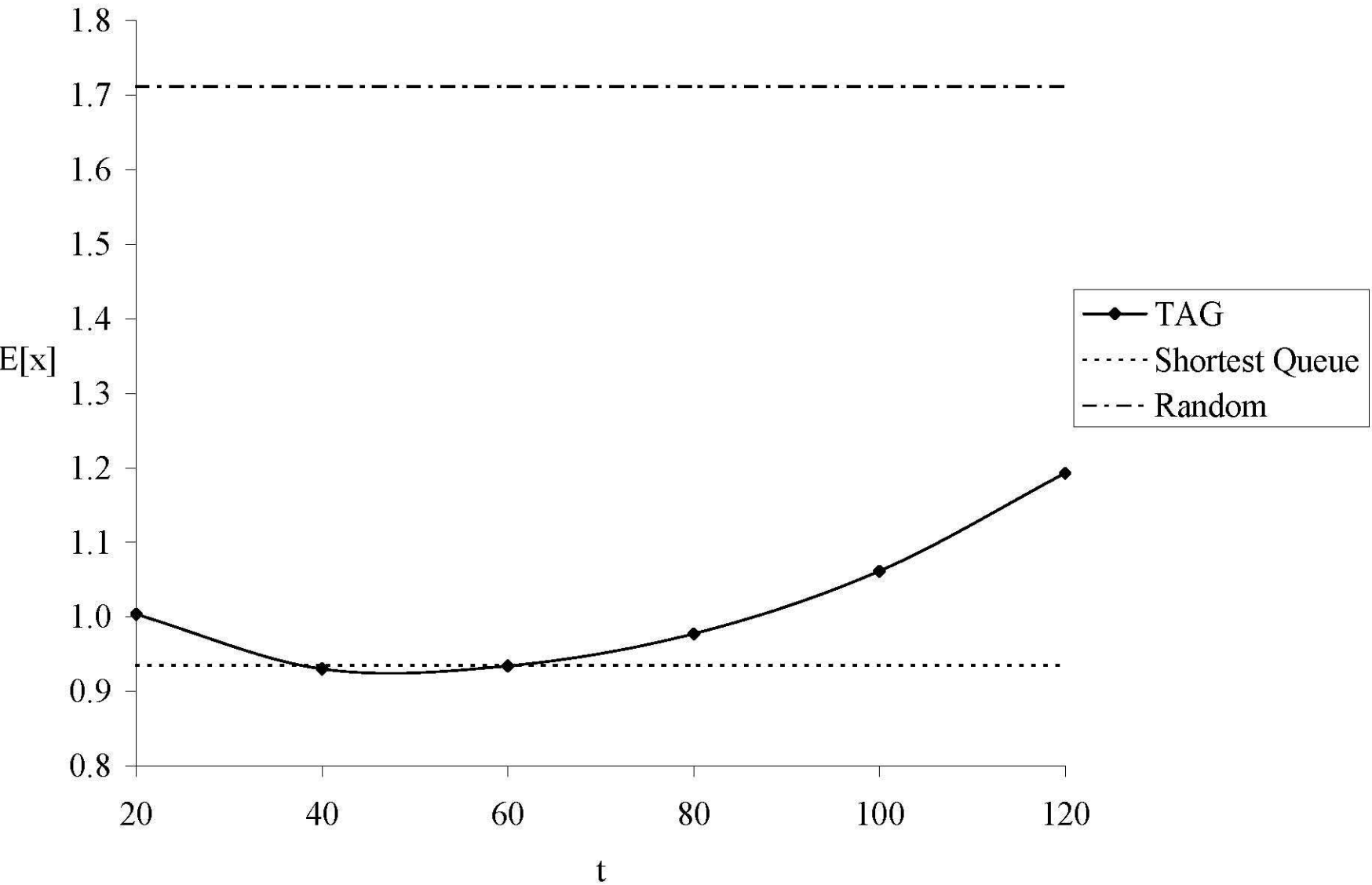
$$Server2 \stackrel{def}{=} (branch, p \times fast).(service2, s_2).Server2 + (branch, (1 - p) \times fast).Server2'$$

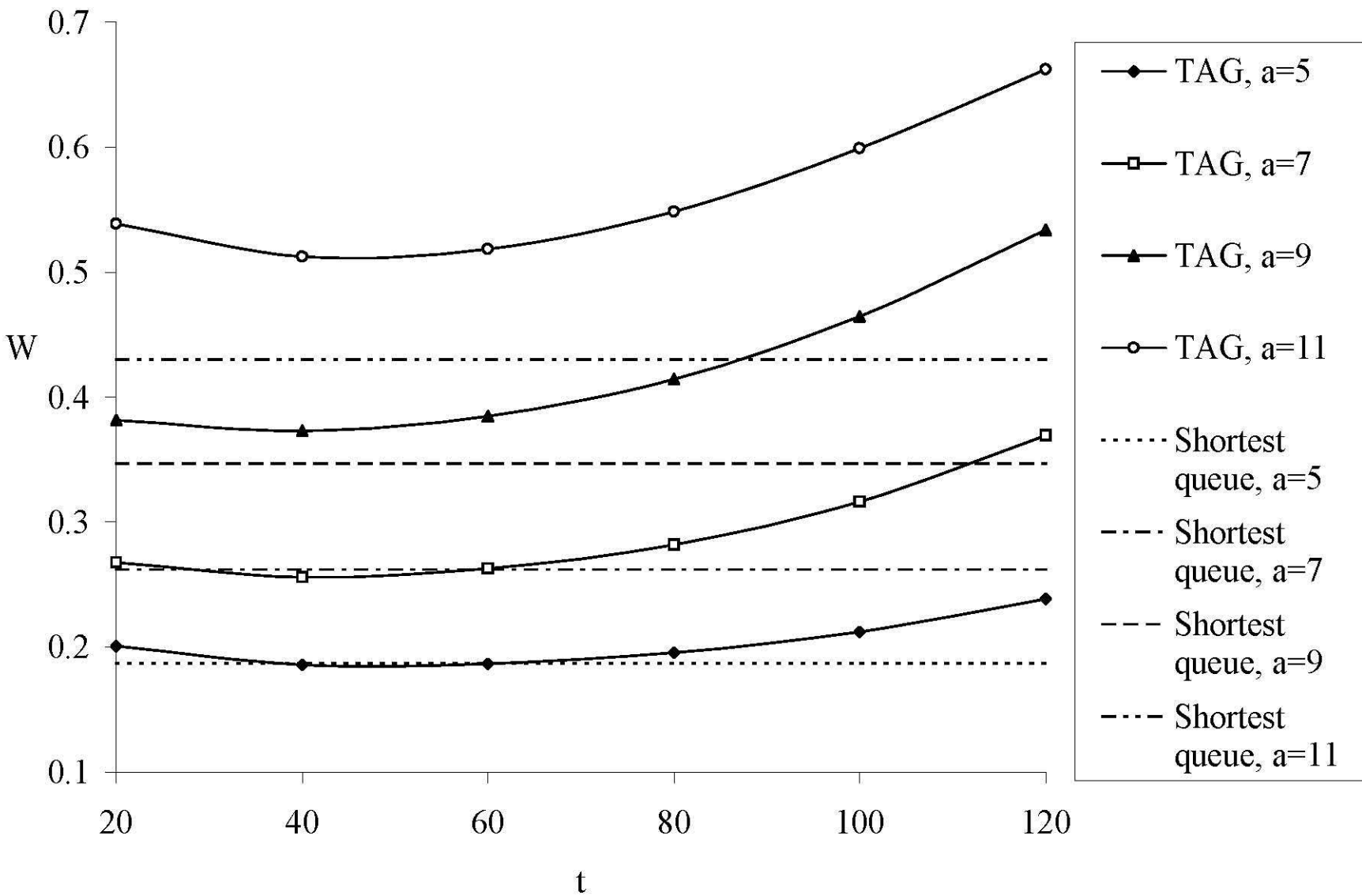
$$Server2' \stackrel{def}{=} (service2, \alpha s_1).Server2 + (phaseone, (1 - \alpha)s_1).(service2, s_2).Server2$$

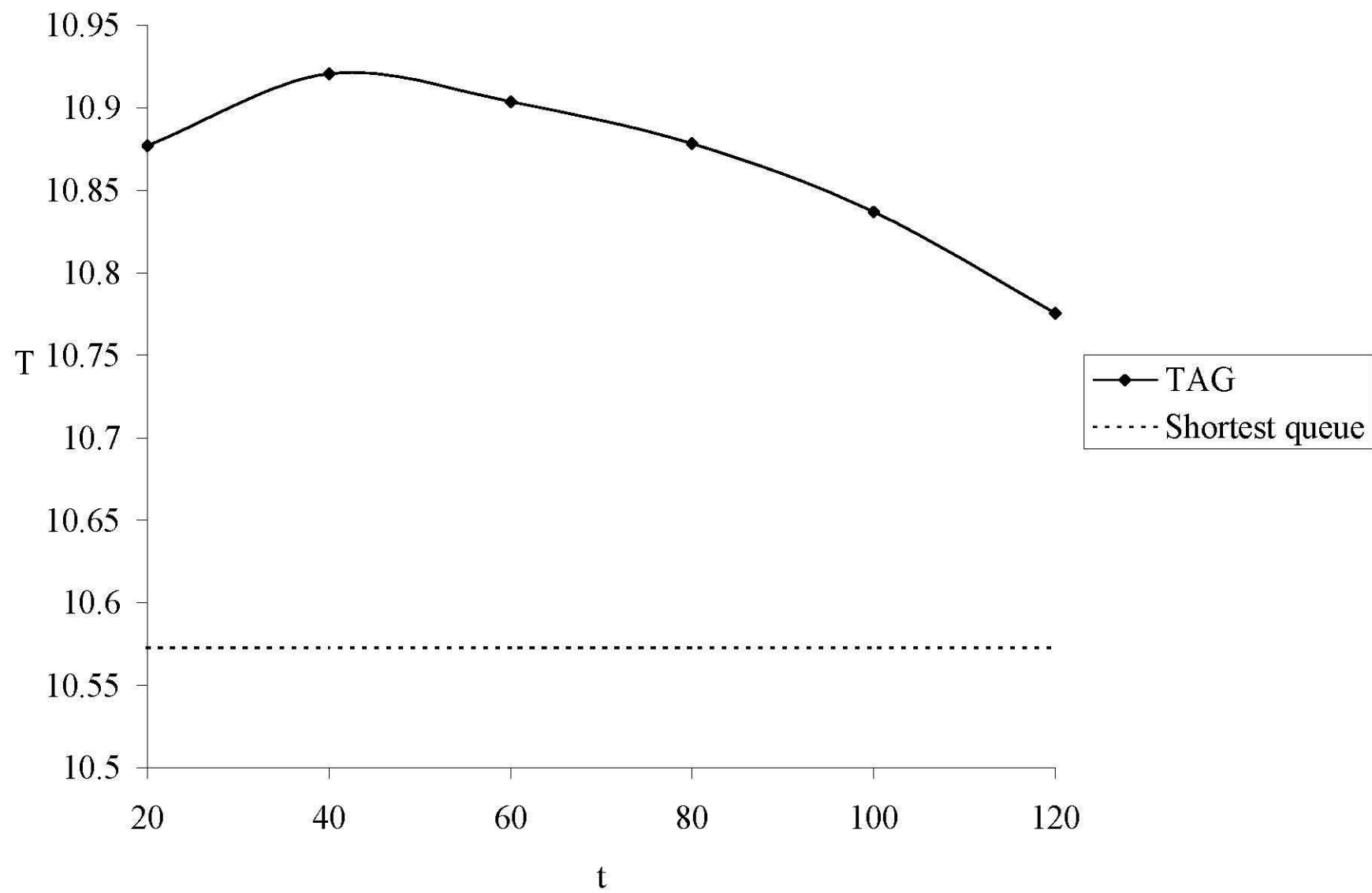
- Actions *phase1* and *branch* must be added to the cooperation sets and queue components to prevent them firing when the queue is empty.
- The *timeout* action is used to reset the service at node 1 if it is interrupted.
- p is the result of the race between *timeout* and *phase1* conditional on *service1* having not occurred.



More results (1/3)









Conclusions

- This paper has only scratched the surface of this problem, much remains to be done.
- It has been shown that TAG has some potential benefit for finite systems.
- More experimentation is required with the Cox distribution, particularly in varying α .
- A simulator has been produced (by Alex Liu) which will be used for comparison and to assess the degree of error introduced by the approximations.

Thanks for your time!

