

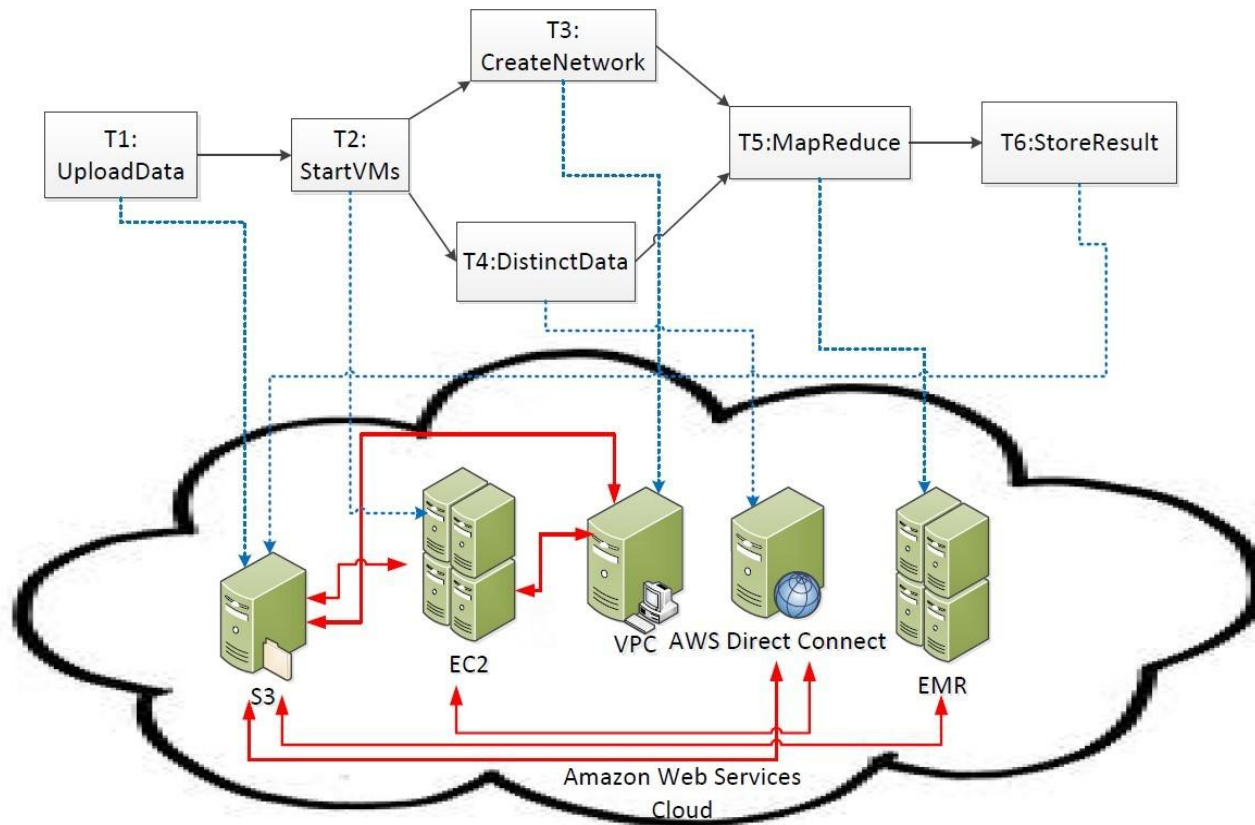
Developing Communication-aware Service Placement Frameworks in the Cloud Economy

Chao Chen, Ligang He, **Bo Gao**, Stephen A. Jarvis

Multitenant cloud system

- Cloud system with multiple tenants and services
 - Cloud Tenants: Users renting services and virtual machines.
 - Cloud system(Datacenter)
 - Cloud providers: deliver a level of QoS(Quality-of-Service), such as: Amazon Web Services, Microsoft Azure, Google Cloud Platform, etc,
 - Cloud Services: hosted by a collection of virtual machines running different jobs: data analysis, web servers, etc.

Example: NASDAQ OMX



S3:
Amazon simple Storage Service

EC2:
Amazon Elastic Compute

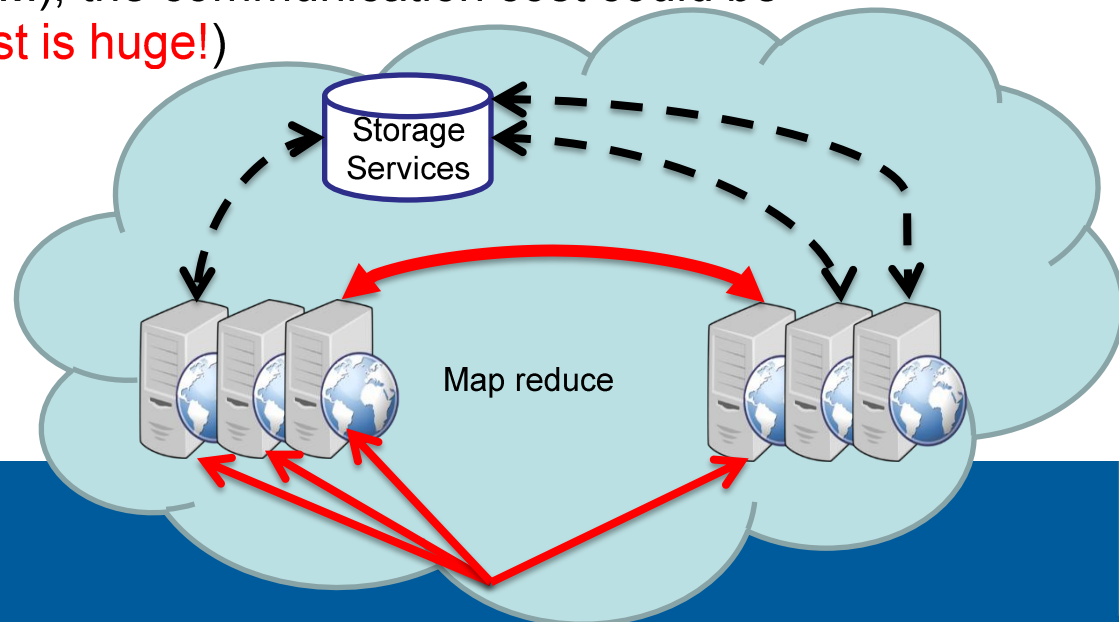
VPC:
Amazon Virtual Private Cloud

Direct Connect:
Amazon Direct Connect

EMR:
Amazon Elastic MapReduce

Challenges

1. Task/Service invocations may vary according to dynamic system information, and it may be difficult to know the full picture of the tasks/workflows in the Cloud. This work is *service-oriented*, which does not focus on allocating resources for a set of specific tasks or workflows, but aim to allocate resources based on *the interaction patterns between services*.
2. When the services interact with each other, data might be(cached) communicated between them. If the Virtual machine(VMs) that host the services with frequent communications among themselves can be placed to the same Physical Machine(PM), the communication cost could be significantly reduced. (**The cost is huge!**)



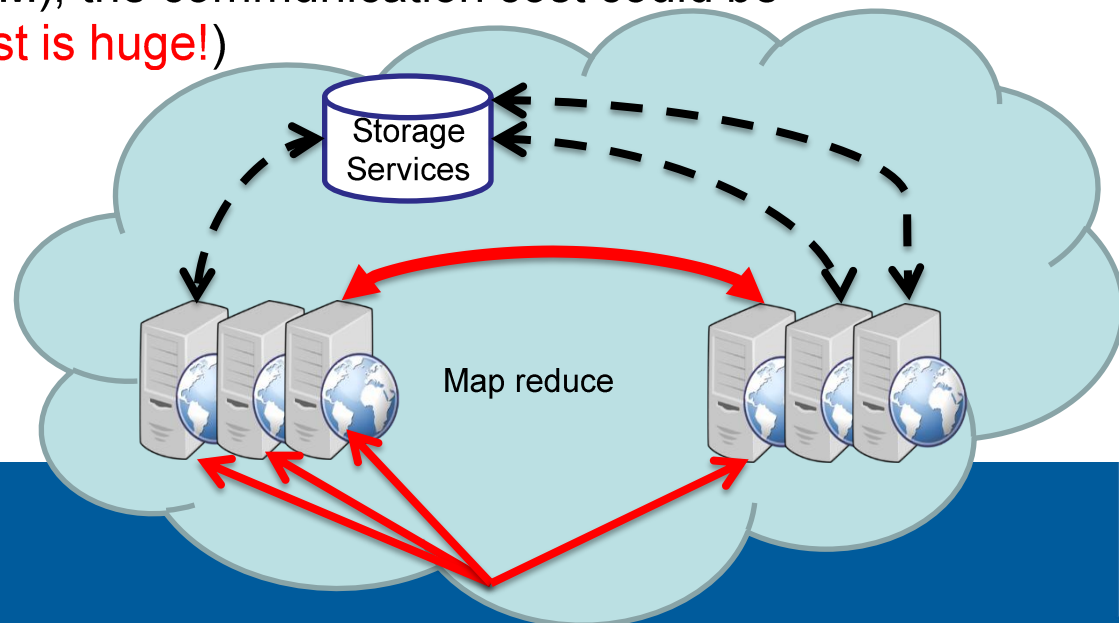
Challenges

1. Task/Service invocations may vary according to dynamic system information, and it may be difficult to know the full picture of the tasks/workflows in the Cloud. This work is *service-oriented*, which does not focus on allocating resources for a set of specific tasks or workflows, but aim to allocate resources based on *the interaction patterns between services*.

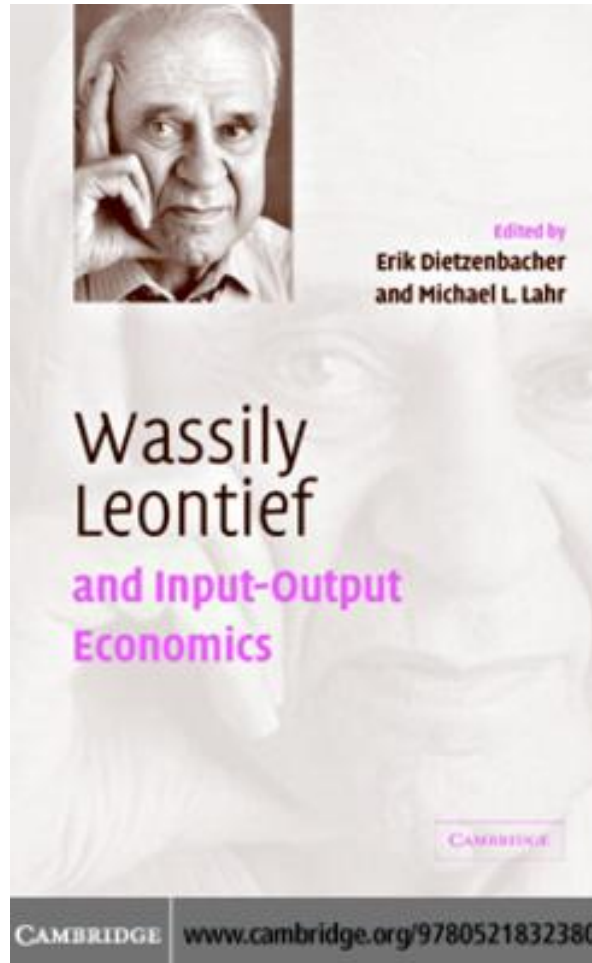
Leontief Input-Output Model in economy

2. When the services interact with each other, data might be(cached) communicated between them. If the Virtual machine(VMs) that host the services with frequent communications among themselves can be placed to the same Physical Machine(PM), the communication cost could be significantly reduced. (**The cost is huge!**)

Genetic algorithm to find out the “optimal” solution



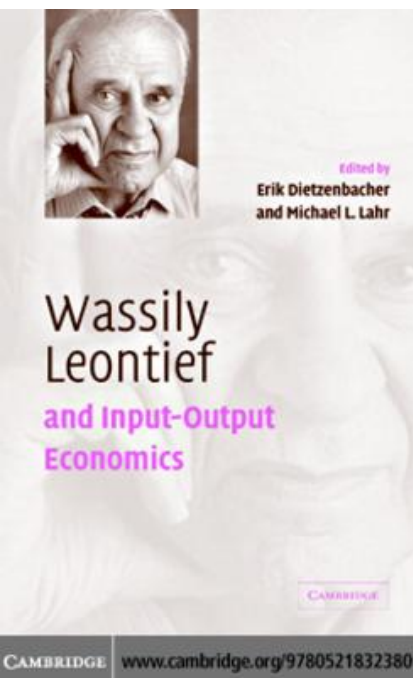
Leontief Input-Output Model



Leontief Input-Output Model

Exchange of Goods and Services in the U.S. for 1947 (in billions of 1947 dollars)

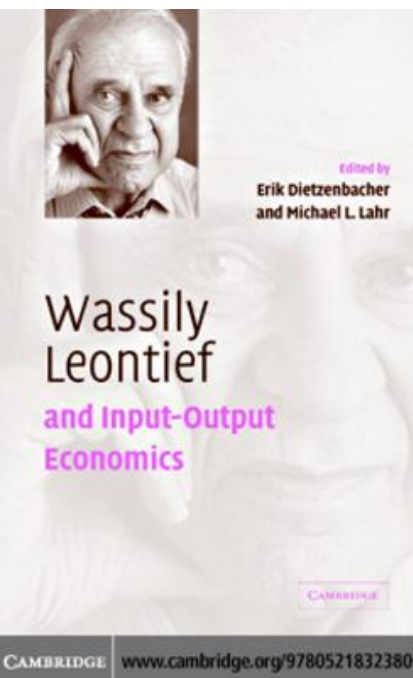
	Agriculture	Manufacturing	Services	Open Sector
Agriculture	34.69	4.92	5.62	39.24
Manufacturing	5.28	61.82	22.99	60.02
Services	10.45	25.95	42.03	130.65
Total Gross Output	84.56	163.43	219.03	



Leontief Input-Output Model

Exchange of Goods and Services in the U.S. for 1947 (in billions of 1947 dollars)

	Agriculture	Manufacturing	Services	Open Sector
Agriculture	34.69	4.92	5.62	39.24
Manufacturing	5.28	61.82	22.99	60.02
Services	10.45	25.95	42.03	130.65
Total Gross Output	84.56	163.43	219.03	



Leontief Input-Output Model

Exchange of Goods and Services in the U.S. for 1947 (in billions of 1947 dollars)

	Agriculture	Manufacturing	Services	Open Sector
Agriculture	34.69	4.92	5.62	39.24
Manufacturing	5.28	61.82	22.99	60.02
Services	10.45	25.95	42.03	130.65
Total Gross Output	84.56	163.43	219.03	

consumption matrix

$$C = \begin{bmatrix} .4102 & .0301 & .0257 \\ .0624 & .3783 & .1050 \\ .1236 & .1588 & .1919 \end{bmatrix}$$

demand vector

$$d = \begin{bmatrix} 39.24 \\ 60.02 \\ 130.65 \end{bmatrix}$$



edited by
Erik Dietzenbacher
and Michael L. Lahr

Wassily
Leontief
and Input-Output
Economics

CAMBRIDGE

CAMBRIDGE | www.cambridge.org/9780521832380

Leontief Input-Output Model

Exchange of Goods and Services in the U.S. for 1947 (in billions of 1947 dollars)

	Agriculture	Manufacturing	Services	Open Sector
Agriculture	34.69	4.92	5.62	39.24
Manufacturing	5.28	61.82	22.99	60.02
Services	10.45	25.95	42.03	130.65
Total Gross Output	84.56	163.43	219.03	

consumption matrix

$$C = \begin{bmatrix} .4102 & .0301 & .0257 \\ .0624 & .3783 & .1050 \\ .1236 & .1588 & .1919 \end{bmatrix}$$

demand vector

$$\mathbf{d} = \begin{bmatrix} 39.24 \\ 60.02 \\ 130.65 \end{bmatrix}$$

\mathbf{x} - Equilibrium Production Level

$$\mathbf{x} = C\mathbf{x} + \mathbf{d}$$

$$\mathbf{x} = (I - C)^{-1}\mathbf{d} = \begin{bmatrix} 82.40 \\ 138.85 \\ 201.57 \end{bmatrix}$$



edited by
Erik Dietzenbacher
and Michael L. Lahr

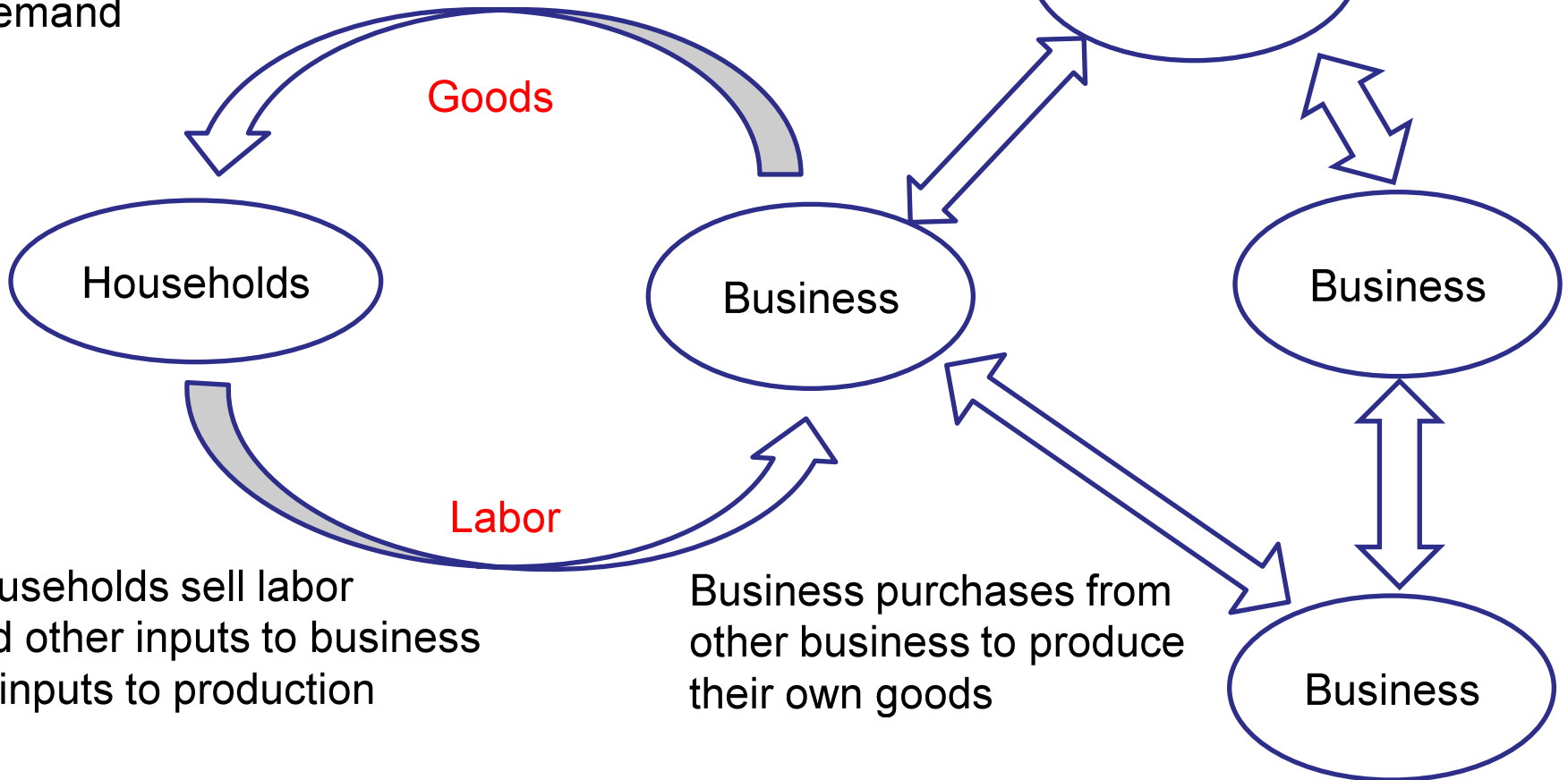
Wassily
Leontief
and Input-Output
Economics

CAMBRIDGE

CAMBRIDGE | www.cambridge.org/9780521832380

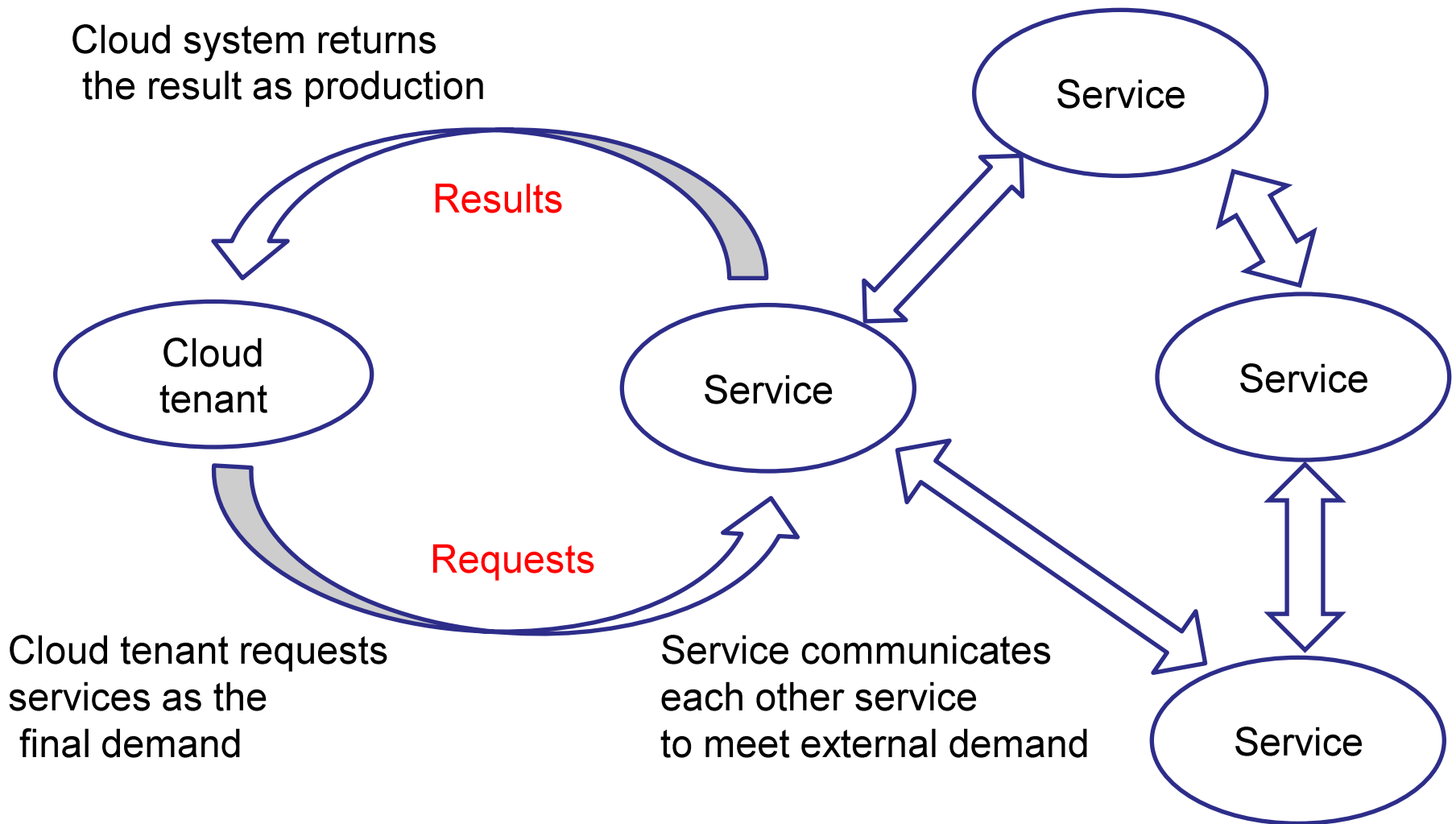
Modern industry ecosystem

Households buy output of business as the final demand



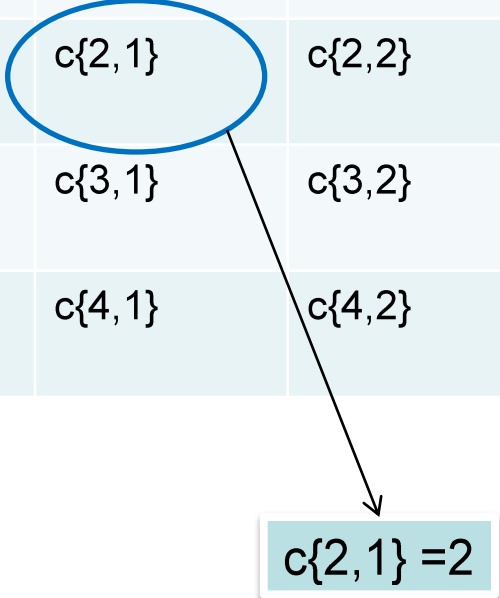
Households sell labor and other inputs to business as inputs to production

Business purchases from other business to produce their own goods

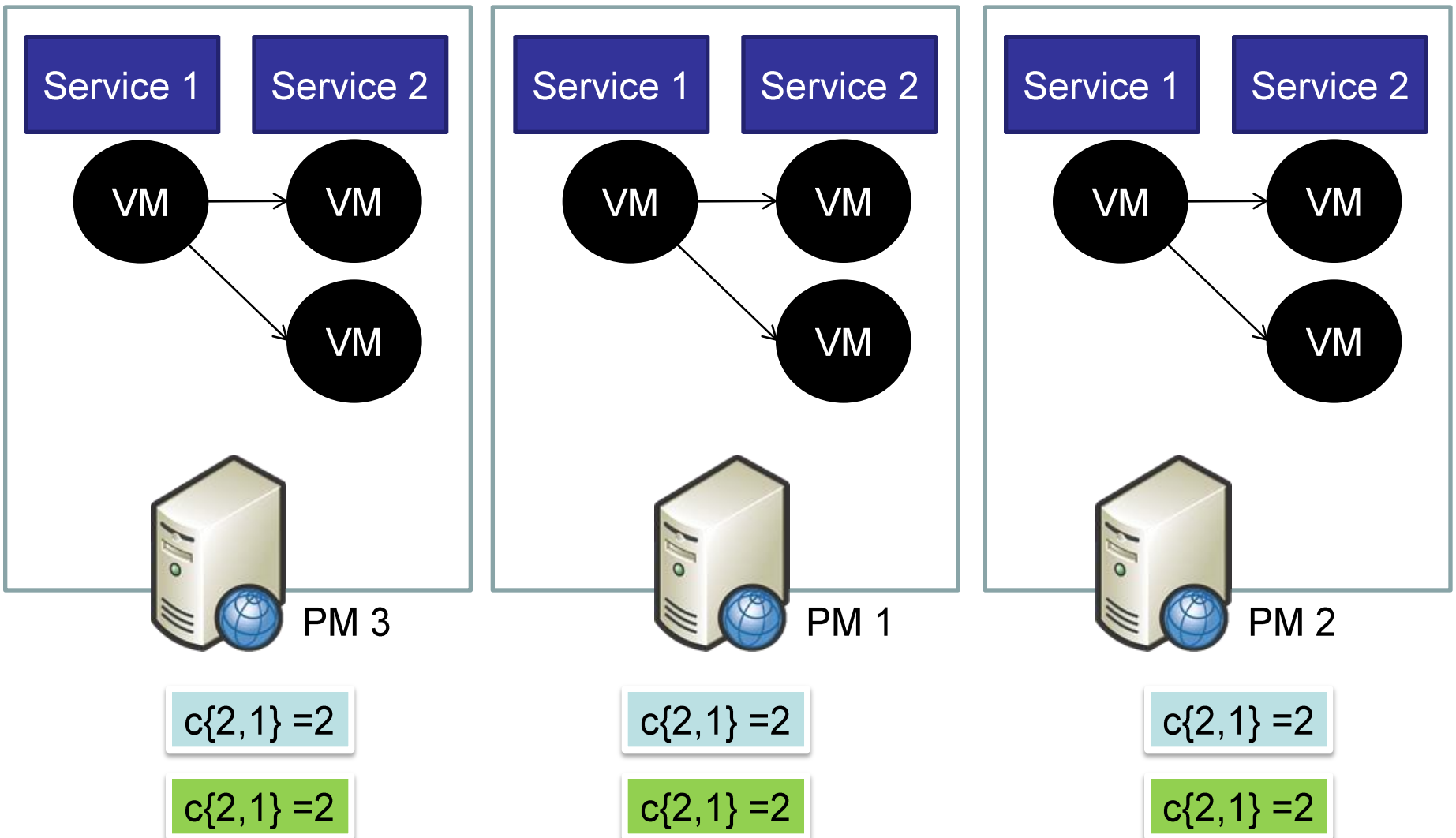


Consumption Matrix

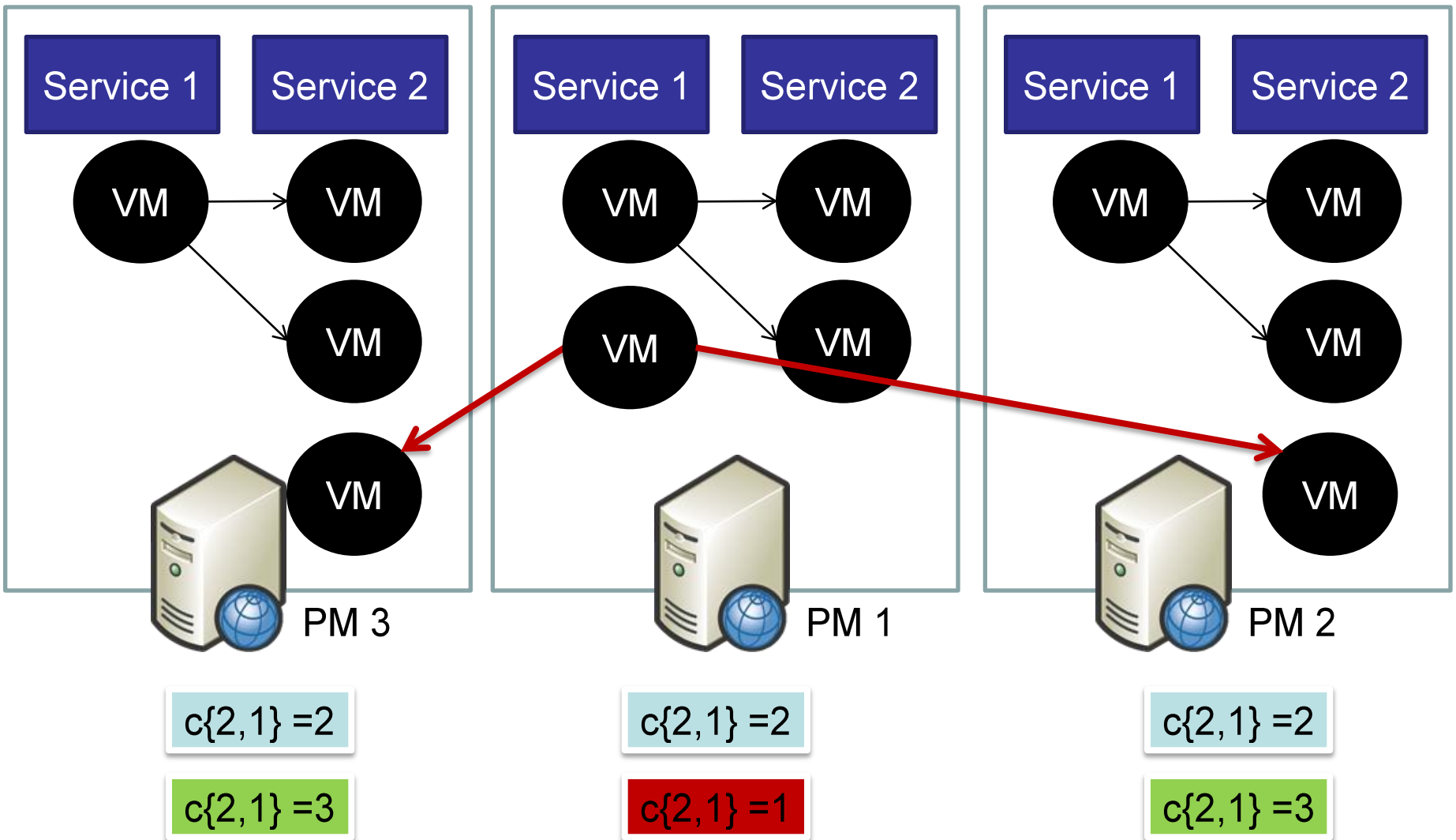
	Service 0	Service 1	Service 2	Service 3	Service 4	External demand
Service 0	$c\{0,0\}$	$c\{0,1\}$	$c\{0,2\}$	$c\{0,3\}$	$c\{0,4\}$	λ_0
Service 1	$c\{1,0\}$	$c\{1,1\}$	$c\{1,2\}$	$c\{1,3\}$	$c\{1,4\}$	λ_1
Service 2	$c\{2,0\}$	$c\{2,1\}$	$c\{2,2\}$	$c\{2,3\}$	$c\{2,4\}$	λ_2
Service 3	$c\{3,0\}$	$c\{3,1\}$	$c\{3,2\}$	$c\{3,3\}$	$c\{3,4\}$	λ_3
Service 4	$c\{4,0\}$	$c\{4,1\}$	$c\{4,2\}$	$c\{4,3\}$	$c\{4,4\}$	λ_4


$$c\{2,1\} = 2$$

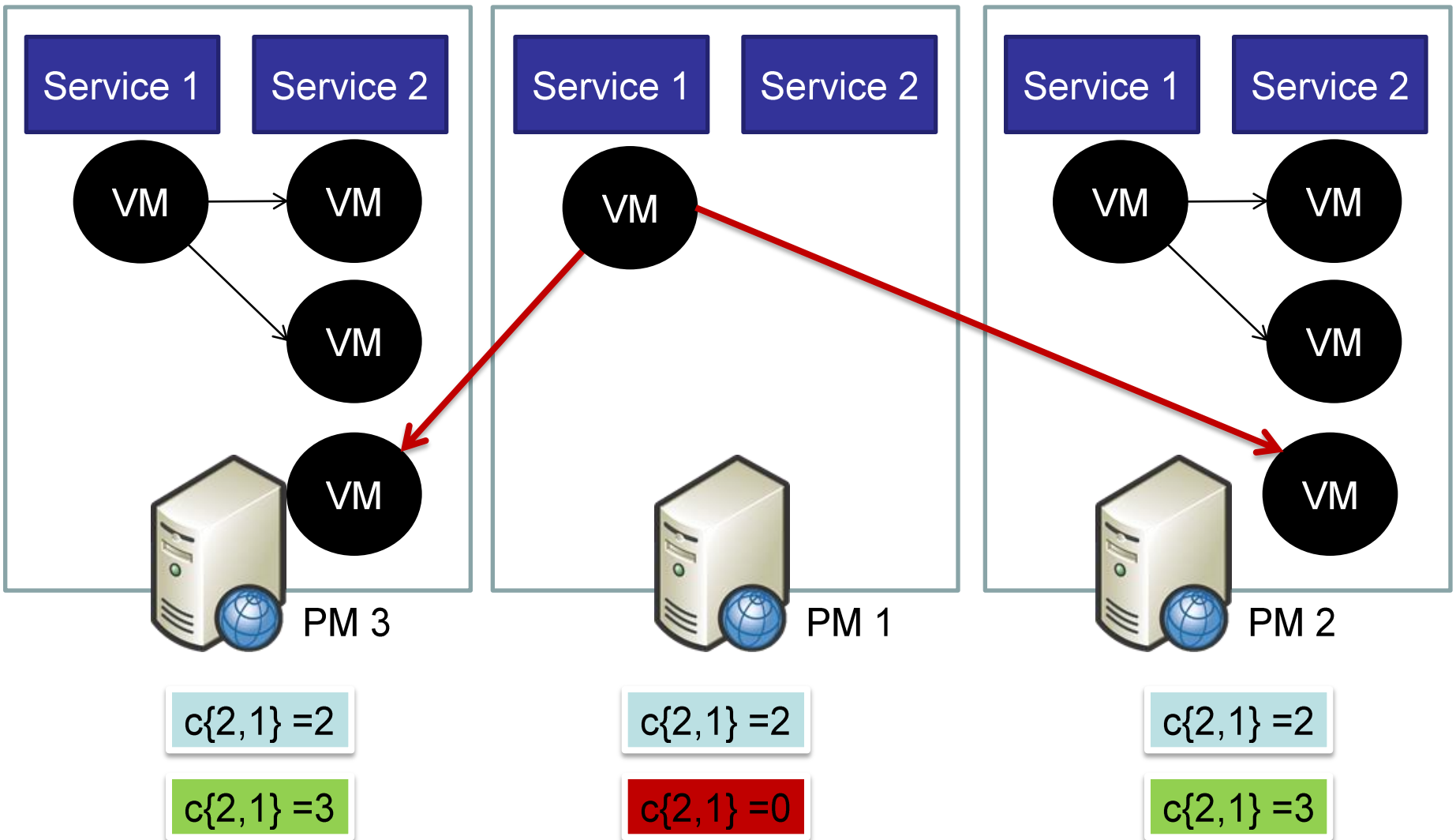
Consumption relation



Consumption relation



Consumption relation



Communication Cost

$$C(\mathcal{M}) = \sum_{k=1}^N \sum_{j=1}^M \sum_{i=1}^M \beta_{ijk} \quad (6)$$

calculates that the amount of requests that are sent from $S\{j\}$ in PM $n\{k\}$ to $S\{i\}$ in a time unit, but cannot be handled by VM $\{i\}$ in $n\{k\}$ in order to maintain the QoS. Therefore, these requests have to be sent to be processed by VM $\{i\}$ in a different PM.

$$\beta_{ijk} = \begin{cases} e_{ji} \times (f(j, R_j, v_{jk}) \times p_{ji} - f(i, R_i, v_{ik})) & \text{if } \alpha_{ijk} < c_{ij} \\ 0 & \text{otherwise} \end{cases} \quad (7)$$

The objective is to find a VM-to-PM mapping such that $C(\mathcal{M})$ is minimized, subject to certain constraints. This can be formalized as Eq. 8, where x_i is the number of VMⁱ's

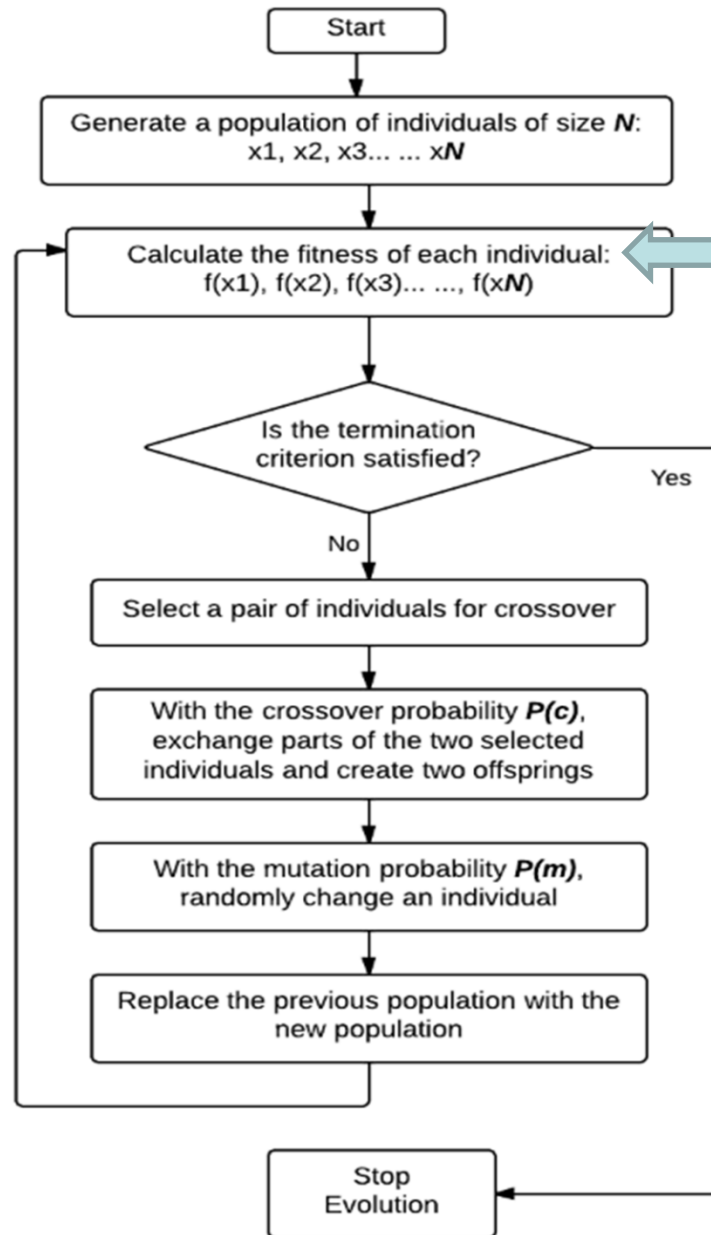
The ratio of the number of VM $\{i\}$ to the number of VM $\{j\}$ in the PM $n\{k\}$

The total amount of data that have to be communicated in the Cloud caused by the inadequate resource capacity of $S\{j\}$ in PM $n\{k\}$ comparing with that of $S\{j\}$ in the same PM.

$$\begin{aligned} & \text{minimize } C(\mathcal{M}), \\ & \text{subject to: } \forall i: 1 \leq i \leq M, \sum_{k=1}^N v_{ik} = x_i \quad (8) \\ & \quad \quad \quad v_{ik} \geq 0 \end{aligned}$$



Designing Genetic Algorithm for VMs Allocation Problem

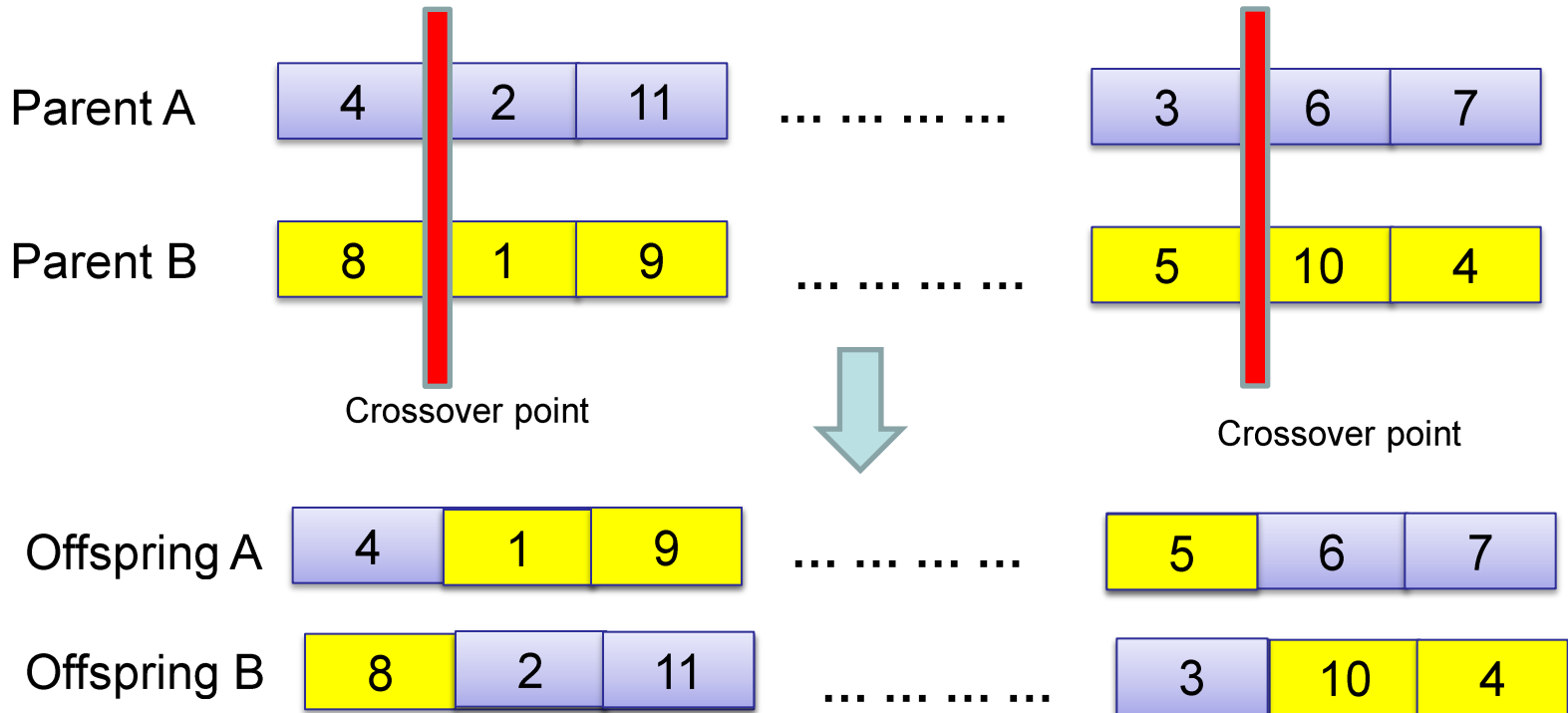


Calculate it via Eq. 5

$$C(\mathcal{M}) = \sum_{k=1}^N \sum_{j=1}^M \sum_{i=1}^M \beta_{ijk}$$

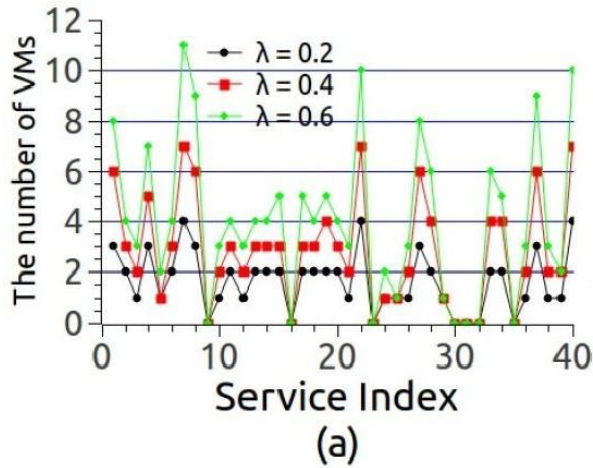


Genetic Algorithm : two points Crossover

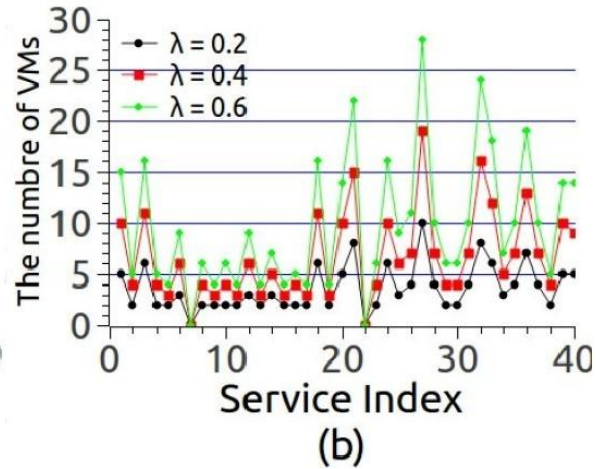


Impact of the increase in external demands

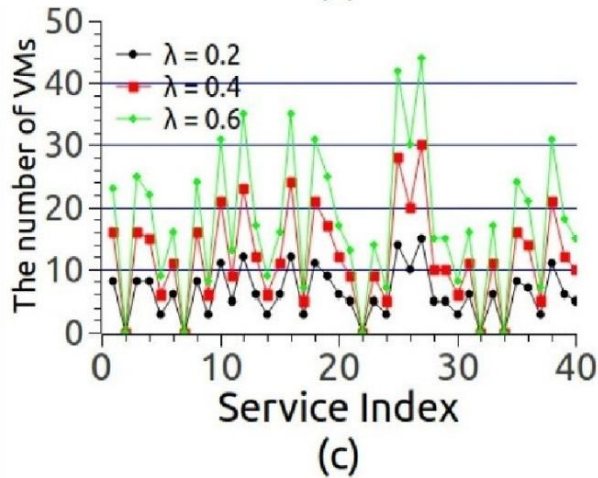
a: computation-intensive workflow



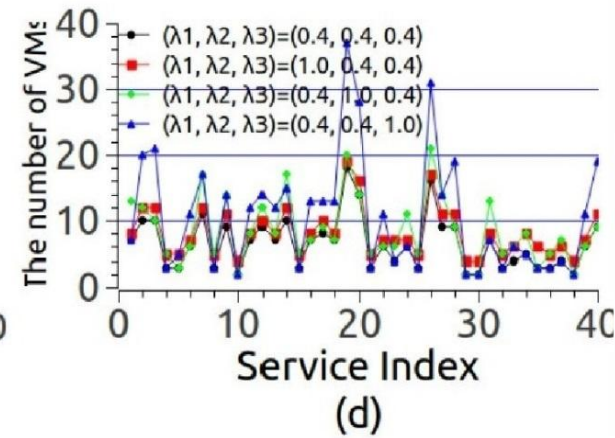
b: general workflow



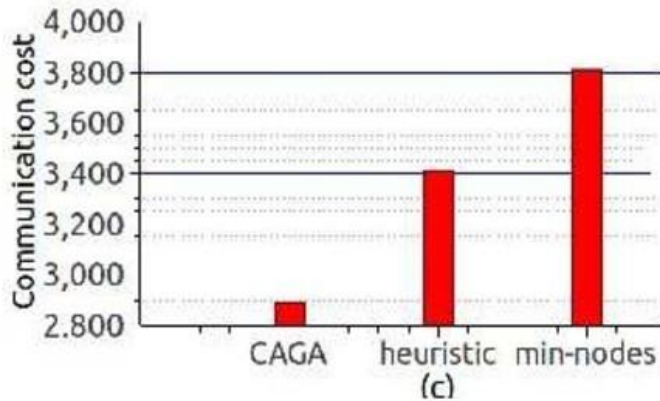
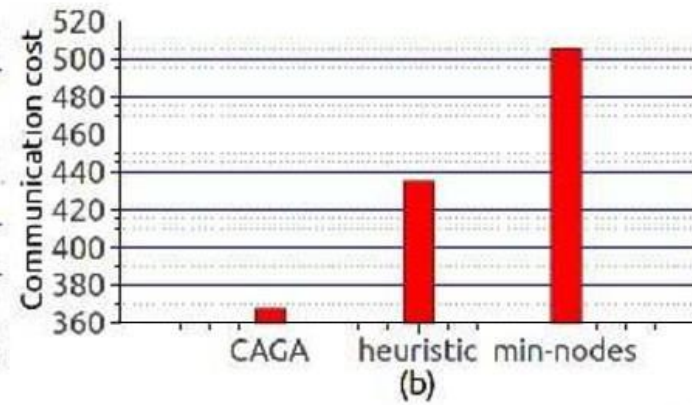
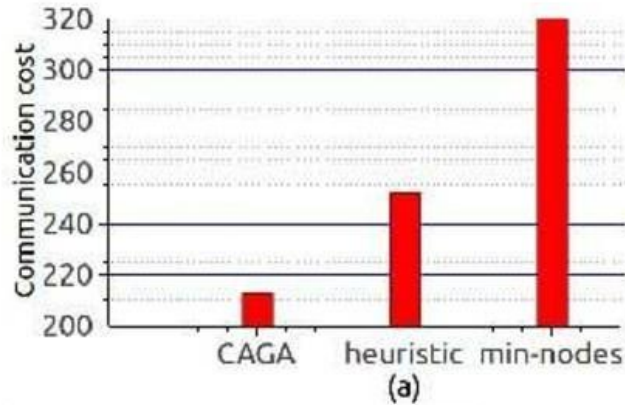
c: communication-intensive workflow



d: the three types of workflow combined



Comparing CAGA with mini-nodes and the round-robin heuristic in terms of communication cost

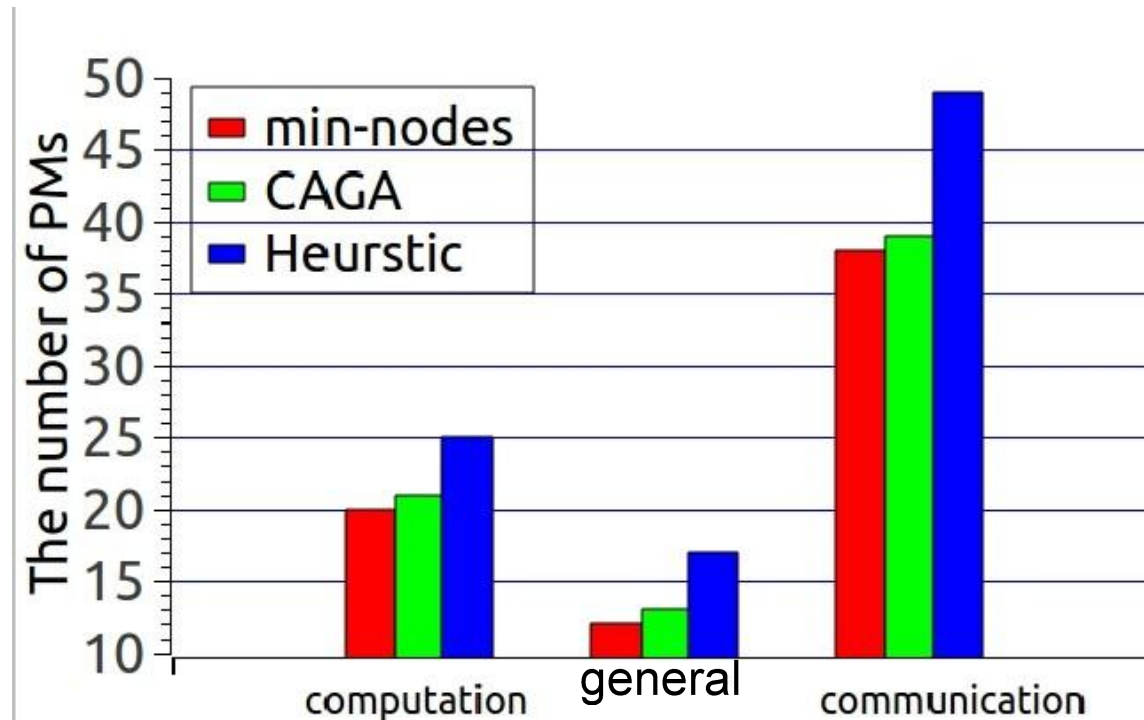


a: computation-intensive workflow

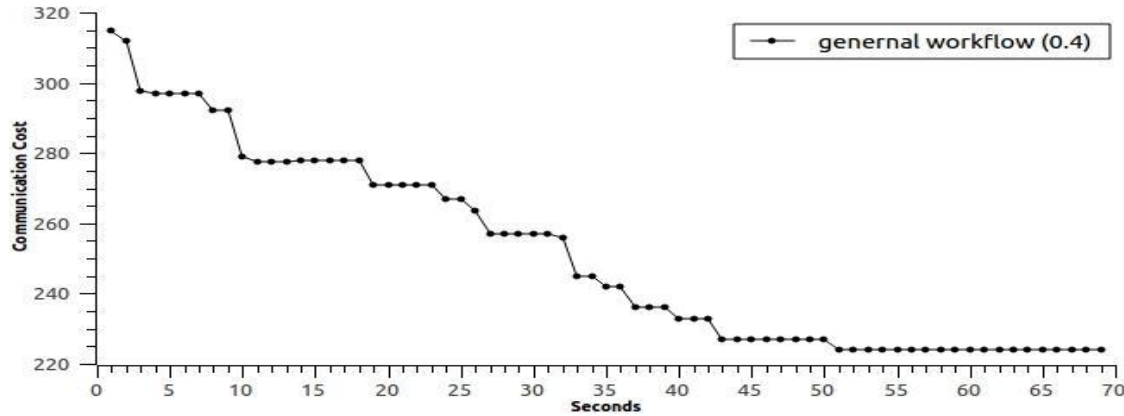
b: general workflow

c: communication-intensive workflow

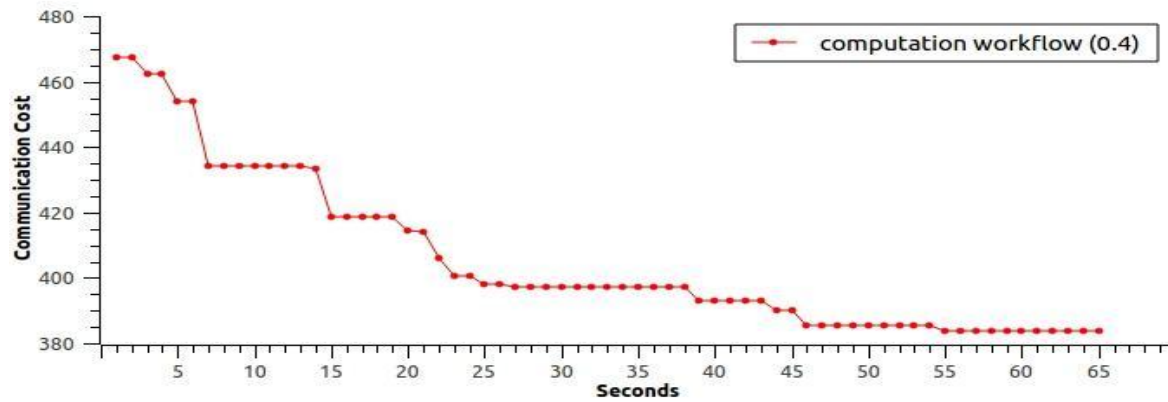
Comparing CAGA with mini-nodes and the round-robin heuristic in terms of number of used PMs



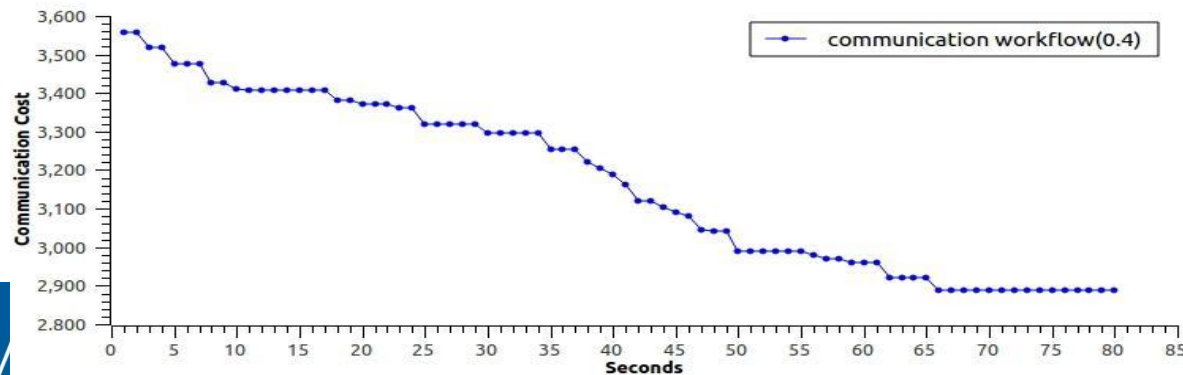
Convergence speed of CAGA



a: general workflow



b: computation-intensive workflow



c: communication-intensive workflow

THANK
YOU!

THE UNIVERSITY OF
WARWICK

Question?

