

Performance Evaluation

Jean-Marc Vincent and Arnaud Legrand ¹

¹Laboratory ID-IMAG
MESCAL Project
Universities of Grenoble
{Jean-Marc.Vincent,Arnaud.Legrand}@imag.fr
<http://www-id.imag.fr/Laboratoire/Membres/>



Outline

- 1 **Scientific context**
- 2 **Methodology**
- 3 **Performance indexes**
- 4 **Experimental framework**
- 5 **Analysis of Experiments**
- 6 **Results synthesis**

Research activities in performance evaluation

Teams in Grenoble

- Mescal project : large systems (clusters and grids)
- Moais project : interactive parallel systems
- Drakkar team : networking
- Sardes : middleware
- Verimag : Embedded systems
- etc

Industrial collaborations

- France-Télécom R & D : load injectors, performances of middlewares
- HP-Labs : cluster computing, benchmarking
- Bull : benchmarking, performances analysis
- ST-Microelectronics

Application context (1)

Complexity of computer systems

- **hierarchy** : level decomposition : OS / Middleware / Application
- **distribution** : asynchronous resources : memory, CPU, network
- **dynamicity** : architecture and environment (reliability, mobility,...)
- **scalability** : number of components (autonomous management)

Typical problems

- Minimize losses in routing policies
- Minimize active waiting in threads scheduling
- Maximize cache hits
- Optimise block sizes in parallel applications
- Maximize throughput of communication systems
- Fix time-outs, reemission periods, ...
- Fix the granularity : pages, blocks, tables, message sizes...
- ...

Application context (2)

Typical “hot” applications

- **Peer to peer systems** : dimensionning, control
- **Mobile networks** : ad-hoc networking, reactivity, coherence
- **Grids** : resources utilization, scheduling
- etc

Other application domains

- production systems : production lines, logistic,...
- embedded systems
- modelling of complex systems : biology, sociology,...
- etc

Development of parallel/distributed applications

- **Qualitative specifications** : Is the result correct ?
 - properties verifications : formal/automatic proofs
 - testing : critical dataset
- **Quantitative specifications** : Is the result obtained in an acceptable time ?
 - performance model
 - performance measurements
- **Problem identification**
 - debugging, log analysis
 - performance statistical analysis
- **Modification**
 - source code / libraries / OS / architecture
 - parameters of the system : dimensioning
 - control algorithms : tuning

Dual analysis

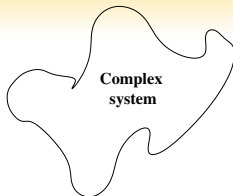
Understand the behavior of a distributed application

- 1 identification of distributed patterns, states of the system
- 2 pattern verification
- 3 time evaluation
- 4 global analysis of the execution and performance synthesis
- 5 system monitoring
- 6 **global cost evaluation for the application user**

Understand resources utilization

- 1 hierarchical model of resources
- 2 evaluation of utilization at :
application level; executive runtime;
operating system; hardware architecture
- 3 **global cost evaluation for the resources manager**

Methodology (1)



Evaluation methods

From abstraction to physical reality

Model

Method

Model

Method

Mathematical \longrightarrow

Analysis (formal, numerical, approximation)

Steps for a Performance Evaluation Study (Jain)

- 1 State the goals of the study : level of decision, investment, optimization, technical,...
- 2 Define system boundaries.
- 3 List system services and possible outcomes.
- 4 Select performance metrics.
- 5 List system and workload parameters
- 6 Select factors and their values.
- 7 Select evaluation techniques.
- 8 Select the workload.
- 9 Design the experiments.
- 10 Analyze and interpret the data.
- 11 Present the results. Start over, if necessary.

Aim of the course

Objective

- 1 Be able to analyze and predict performances of parallel/distributed systems
- 2 Be able to build a software environment that produces the performances indexes.

Methods

- 1 Specification and identification of problems : modelling
- 2 Analysis of quantitative models : formal, numerical, simulation
- 3 Experimentation and statistical data analysis.

Organization of the course

Practical evaluation of systems 5 lectures 3h

- 1 Friday 14/10/2011 (13h30-16h45): Jean-Marc Vincent] Introduction, performance indexes, data analysis, modeling and inference.
- 2 Friday 21/10/2011 (13h30-16h45): Arnaud Legrand] Measurement on computer systems (benchmarking, observation, tracing, monitoring, profiling).
- 3 Friday 4/11/2011 (13h30-16h45): Arnaud Legrand] Visualization and discrete event simulation of computer systems.
- 4 Monday 7/11/2011 (13h30-16h45): Arnaud Legrand and Jean-Marc Vincent] Emulation of computer systems and random number generation.
- 5 Friday 18/11/2011 (13h30-16h45): Jean-Marc Vincent and Arnaud Legrand] Workload generation and introduction to design of experiments.

Evaluation

Reading of an article, synthesis and presentation

References : text books

- **The Art of Computer Systems Performance Analysis : Techniques for Experimental Design, Measurement, Simulation and Modeling.** Raj Jain *Wiley 1991 (nouvelles versions)*
Covers the content of the course, a complete book
- **Performance Evaluation** Jean-Yves Le Boudec EPFL electronic book
<http://ica1www.epfl.ch/perfeval/lectureNotes.htm>
Covers the statistical part of the course
- **Measuring Computer Performance: A Practitioner's Guide** David J. Lilja *Cambridge University press 2000*
Covers the practical part of measurement and benchmarking
- **Discrete-Event System Simulation** Jerry Banks, John Carson, Barry L. Nelson, David Nicol, *Prentice Hall, 2004*
Covers the part on simulation

References : journals and conferences

- **General:** JACM, ACM Comp. Surv., JOR, IEEE TSE,...
- **Specialized:** Performance Evaluation, Operation research, MOR, ACM TOMACS, Queueing Systems, DEDS, ...
- **Application:** IEEE TPDS, TC, TN, TAC, Networks,...
- **Theoretical:** Annals of Probability, of Appl. Prob, JAP, Adv. Appl. Prob,...
- **Conferences on performances:** Performance, ACM-SIGMETRICS, TOOLS, MASCOT, INFORMS, ...
- **Conferences on an application domain:** ITC, Europar, IPDPS, Renpar, ...
- **National seminars:** Atelier d'évaluation de performances,...

Networking

Flow performance

- latency, waiting time, response time
- loss probability
- jitter

Operator performance

- bandwidth utilisation
- achievable throughput
- loss rate

Quality of service

contract between user and provider

service guarantees

tradeoff between utilization and QoS

Parallel processing

Program execution

- makespan, critical path
- speedup, efficiency
- active waiting, communication overlapping
- throughput

System utilization

- cpu utilization, idle time
- memory occupancy
- communication throughput

Parallel programming and scheduling

granularity of the application

tradeoff between utilization and makespan

Distributed applications

Application

- response time
- reactivity
- throughput (number of processed requests/unit time)
- streaming rate

System utilization

- service availability
- resource utilization
- communication throughput

System security

- reliability (error-free period)
- availability

Synthesis

User point of view

optimize its own performance

- get the maximum amount of resources for its own purpose
- guarantee the higher quality of service

Resource point of view

Contract between users and resources:

- guarantee of "equity"
- optimize the use of resources
- minimize costs by identifying performance bottlenecks

Tradeoff Performance - Cost

Why experiments ?

Design of architectures, softwares

- System debugging (!!)
- Validation of a proposition
- Qualification of a system
- Dimensioning and tuning
- Comparison of systems

Many purposes \Rightarrow different methodologies

Experiments fundamentals

Scientific Method

Falsifiability is the logical possibility that an assertion can be shown false by an observation or a physical experiment. [Popper 1930]

Modelling comes before experimenting

Modelling principles [J-Y LB]

- (Occam:) if two models explain some observations equally well, the simplest one is preferable
- (Dijkstra:) It is when you cannot remove a single piece that your design is complete.
- (Common Sense:) Use the adequate level of sophistication.

Design of experiments (introduction)

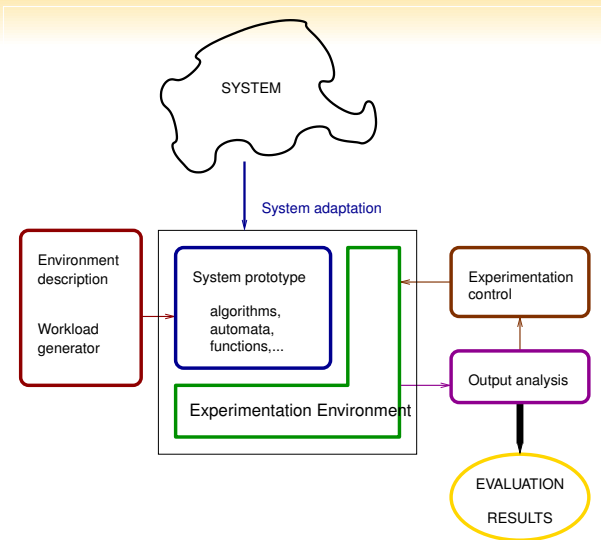
Formulation of the question

Give explicitly the question (specify the context of experimentation)

- Identify parameters (controlled and uncontrolled)
- Identify factors (set levels)
- Specify the response of the experiment

Minimize the number of experiments for a maximum of accuracy

Experimental Framework



Observation technique

Integrated environment : Benchmarks

- Qualification
- Comparison
- Standardization

No interpretation

Level of observation

- Instruction level (Papi)
- System level (OS probes)
- Middleware level (JVMTI)
- Application level (traced libraries, MPITrace)
- User level (own instrumentation point)

Build a semantic on events

Qualification of experiments

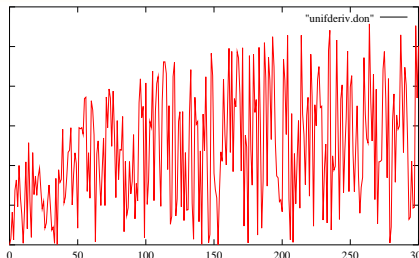
Qualification of measurement tools

- Correctness
- Accuracy
- Fidelity
- Coherence (set of tools)

Qualification on the sequence of experiments

- Reproducibility
- Independence from the environment
- Independence one with each others

Control of experiments (1)



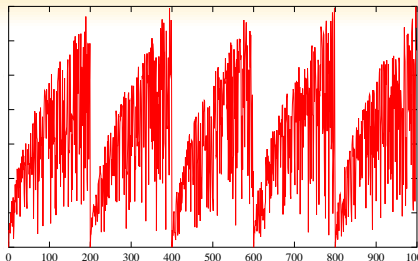
Tendency analysis

non homogeneous experiment

⇒ model the evolution of experiment
estimate and compensate tendency

explain why

Control of experiments (2)



Periodicity analysis

periodic evolution of the experimental environment ?

⇒ model the evolution of experiment

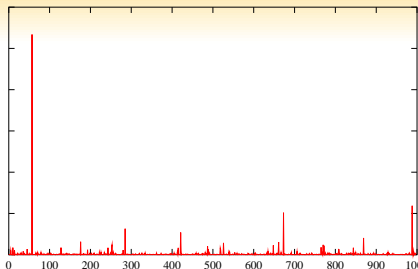
Fourier analysis of the sample

Integration on time (sliding window analysis) Danger : size of the window

Wavelet analysis

explain why

Control of experiments (3)



Non significant values

extraordinary behaviour of experimental environment

rare events with different orders of magnitude

⇒ threshold by value

Danger : choice of the threshold : indicate the rejection rate

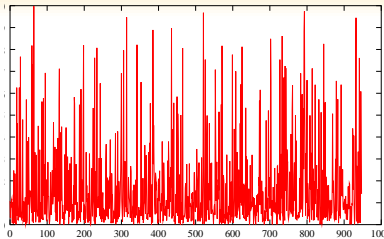
⇒ threshold by quantile

Danger : choice of the percentage : indicate the rejection value

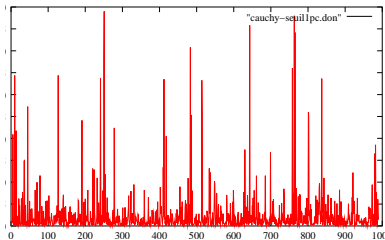
explain why

Control of experiments (4)

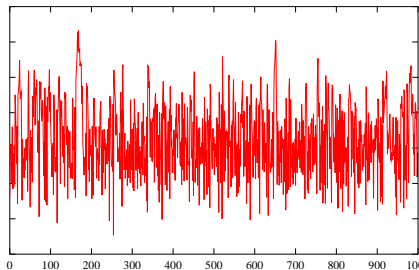
Threshold value : 10



Threshold percentage : 1%



Control of experiments (5)



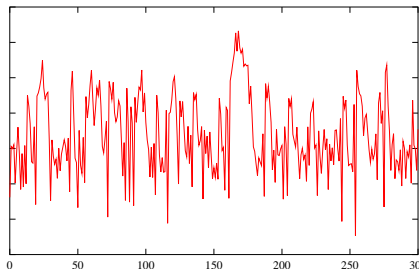
looks like correct experiments

Statistically independent

Statistically homogeneous

Control of experiments (5bis)

Zooming



Autocorrelation

Danger time correlation among samples

experiments impact on experiments

⇒ stationarity analysis

autocorrelation estimation (ARMA)

Experimental results

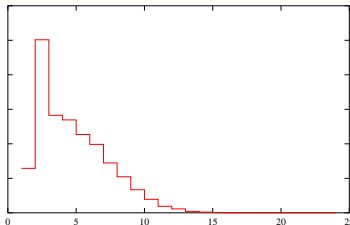
- Deterministic (controlled error non significant (white noise))
- Statistic (the system is non deterministic)

Sample analysis

- Identification of the response set
- Structure of the response set (measure)

Distribution analysis

Summarize data in a **histogram**



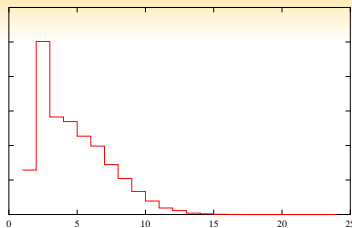
Shape analysis

- unimodal / multimodal
- variability
- symmetric / dissymmetric (skewness)
- flatness (kurtosis)

⇒ **Central tendency analysis**

⇒ **Variability analysis around the central tendency**

Mode value



Mode

- **Categorical data**
- Most frequent value
- highly unstable value
- for continuous value distribution depends on the histogram step
- interpretation depends on the flatness of the histogram

⇒ **Use it carefully**

⇒ **Predictor function**

Median value

Median

- **Ordered data**
- Split the sample in two equal parts

$$\sum_{i \leq \text{Median}} f_i \leq \frac{1}{2} \leq \sum_{i \leq \text{Median}+1} f_i.$$

- more stable value
- does not depends on the histogram step
- difficult to combine (two samples)

⇒ **Randomized algorithms**

Mean value

Mean

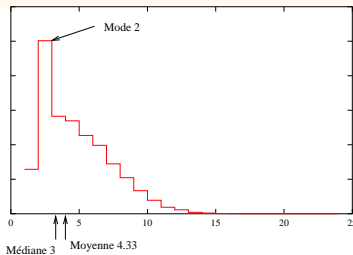
- **Vector space**
- Average of values

$$\text{Mean} = \frac{1}{\text{Sample_Size}} \sum x_i = \sum_x x \cdot f_x.$$

- stable value
- does not depends on the histogram step
- easy to combine (two samples \Rightarrow weighted mean)

\Rightarrow **Additive problems (cost, durations, length,...)**

Central tendency



Complementarity

- Valid if the sample is "Well-formed"
- **Semantic of the observation**
- Goal of analysis

⇒ **Additive problems (cost, durations, length,...)**

Central tendency (2)

Summary of Means

- Avoid means if possible
Loses information
- **Arithmetic mean**
When sum of raw values has physical meaning
Use for summarizing times (not rates)
- **Harmonic mean**
Use for summarizing rates (not times)
- **Geometric mean**
Not useful when time is best measure of perf
Useful when multiplicative effects are in play

Computational aspects

- Mode : computation of the histogram steps, then computation of max $O(n)$ “off-line”
- Median : sort the sample $O(n\log(n))$ or $O(n)$ (subtile algorithm) “off-line”
- Mean : sum values $O(n)$ “on-line” computation

Is the central tendency significant ?
⇒ **Explain variability.**

Variability

Categorical data (finite set)

f_i : empirical frequency of element i

Empirical entropy

$$H(f) = \sum_i f_i \log f_i.$$

Measure the empirical distance with the uniform distribution

- $H(f) \geq 0$
- $H(f) = 0$ iff the observations are reduced to a unique value
- $H(f)$ is maximal for the uniform distribution

Variability (2)

Ordered data

Quantiles : quartiles, deciles, etc

Sort the sample :

$$(x_1, x_2, \dots, x_n) \longrightarrow (x_{(1)}, x_{(2)}, \dots, x_{(n)});$$

$$Q_1 = x_{(n/4)}; \quad Q_2 = x_{(n/2)} = \textit{Median}; \quad Q_3 = x_{(3n/4)}.$$

For deciles

$$d_i = \operatorname{argmax}_i \left\{ \sum_{j \leq i} f_j \leq \frac{i}{10} \right\}.$$

Utilization as quantile/quantile plots to compare distributions

Variability (3)

Vectorial data

Quadratic error for the mean

$$\text{Var}(X) = \frac{1}{n} \sum_1^n (x_i - \bar{x}_n)^2.$$

Properties:

$$\text{Var}(X) \geq 0;$$

$$\text{Var}(X) = \overline{x^2} - (\bar{x})^2, \text{ où } \overline{x^2} = \frac{1}{n} \sum_{i=1}^n x_i^2.$$

$$\text{Var}(X + \text{cste}) = \text{Var}(X);$$

$$\text{Var}(\lambda X) = \lambda^2 \text{Var}(X).$$

A simple example

Maximum value

```
int maximum (int * T, int n)
{ T array of distinct integers,
  {n Size of T}
  {
    int max,i;
    max= int_minimal_value;
    for (i=0; i < n; i++) do
      if (T[i] > max)
        {
          max = T[i];
          Process(max); {Cost of the algorithm}
        }
    end for
    return(max)
  }
```

Cost of the algorithm

Number of calls to **Process**

- minimum : 1
example : $T=[n,1,2,\dots,n-1]$
min cases : $(n-1)!$
- maximum : n
example : $T=[1,2,\dots,n]$
max case : 1

Bounded by a linear function $\mathcal{O}(n)$

But on average ?

A simple example (2)

Theoretical complexity

On average the complexity of the algorithm is :

Build the program

Put probes on the program

Questions :

- 1 Given $n = 1000$ does the observed cost follows the theoretical value ?
- 2 Does the average cost follows the theoretical complexity for all n ?
- 3 Does the average execution time linearly depends on the average cost ?

Modelling

Basic assumptions :

- Data are considered as random variables
- Mutually independent
- Same probability distribution

Check Check Check

The distribution is given by

- Probability density function (pdf) (asymptotic histogram)

$$f_X(x) = \mathbb{P}(x \leq X \leq x + dx) / dx = F'_X(x).$$

- Cumulative distribution function

$$F_X(x) = \mathbb{P}(X \leq x);$$

- Moments : $M_n = \mathbb{E}X^n$, Variance

Average convergence

Law of large numbers

Let $\{X_n\}_{n \in \mathbb{N}}$ be a iid random sequence with finite variance, then

$$\lim_{n \rightarrow +\infty} \frac{1}{n} \sum_{i=1}^n X_i = \mathbb{E}X, \quad \text{almost surely and in } L^1.$$

- convergence of empirical frequencies
- for any experience we get the same result
- fundamental theorem of probability theory

$$\text{Notation : } \bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i.$$

Law of errors

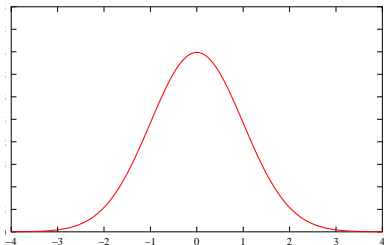
Central limit theorem (CLT)

Let $\{X_n\}_{n \in \mathbb{N}}$ be a iid random sequence with finite variance σ^2 , then

$$\lim_{n \rightarrow +\infty} \frac{\sqrt{n}}{\sigma} (\bar{X}_n - \mathbb{E}X) \stackrel{\mathcal{L}}{=} \mathcal{N}(0, 1).$$

→ error law (Gaussian law, Normal distribution, Bell curve,...)

→ Normalized mean = 0, variance = 1



Distribution

$$\mathbb{P}(X \in [-1, 1]) = 0.68;$$

$$\mathbb{P}(X \in [-2, 2]) = 0.95;$$

$$\mathbb{P}(X \in [-3, 3]) \geq 0.99.$$

Confidence intervals

Confidence level α compute ϕ_α

$$\mathbb{P}(X \in [-\phi_\alpha, \phi_\alpha]) = \alpha$$

For n sufficiently large ($n > 50$)

$$\mathbb{P}\left(\left[\bar{X}_n - \frac{\phi_\alpha \sigma}{\sqrt{n}}, \bar{X}_n + \frac{\phi_\alpha \sigma}{\sqrt{n}}\right] \ni \mathbb{E}X\right) = 1 - \alpha.$$



Confidence intervals (2)

Need an estimator of the variance

$$\hat{\sigma}_{n}^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2.$$

Danger n too small \rightarrow with a normal hypothesis take Student statistic
Three step method

- 1 In a first set of experiments check that the hypothesis is valid
- 2 Estimate roughly the variance
- 3 Estimate the mean and control the number of experiment by a confidence interval

How to report experiments

Problem : provide "nice" pictures to help the understanding

- **Increases deeply the quality of a paper**
- Show the scientific quality of your research
- Observation leads to open problems
- Pictures generates discussions

Mistakes

- **semantic of graphical objects**
- conventions for graphics reading
- first step in scientific validation

Guidelines for good graphics (Jain)

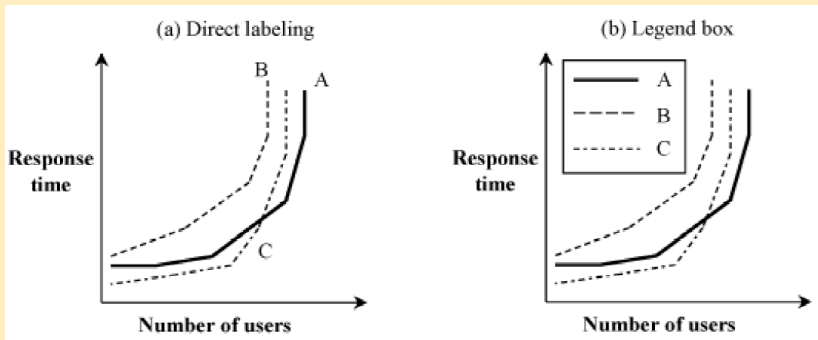
Guidelines for Preparing Good Graphic Charts

Specify the amount of information given by the chart

- 1 Require Minimum Effort from the Reader
- 2 Maximize Information
- 3 Minimize Ink
- 4 Use Commonly Accepted Practices
- 5 Make several trials before arriving at the final chart. Different combinations should be tried and the best one selected.

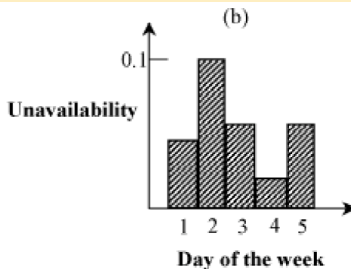
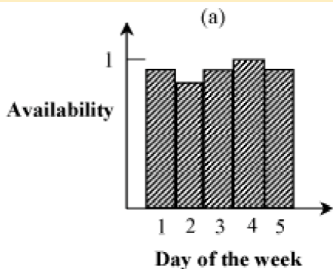
Guidelines for good graphics (Jain)

Minimum effort for the reader



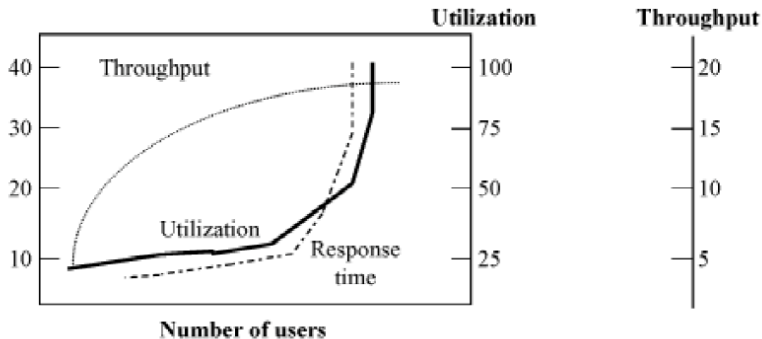
Guidelines for good graphics (Jain)

Minimize Ink



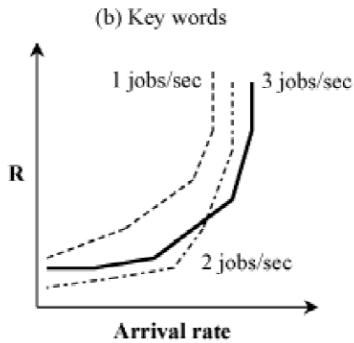
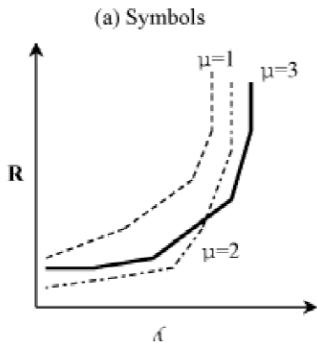
Common mistakes

Multiple scaling, Too much information



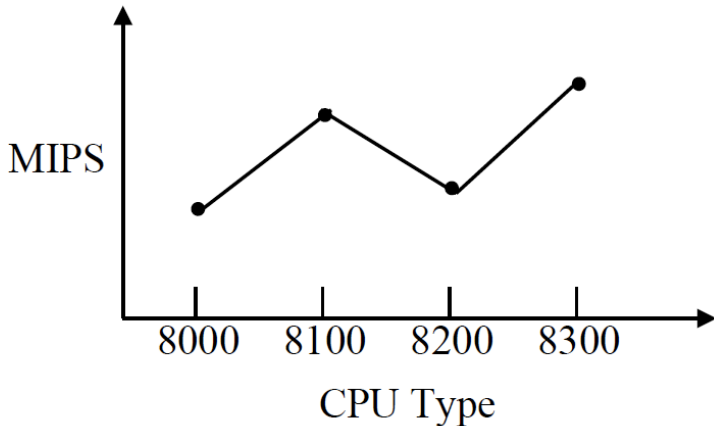
Common mistakes

Cryptic information



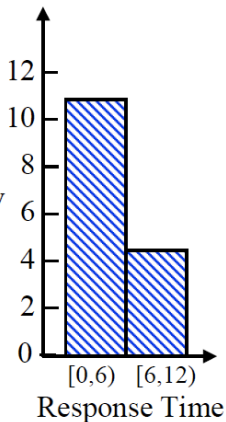
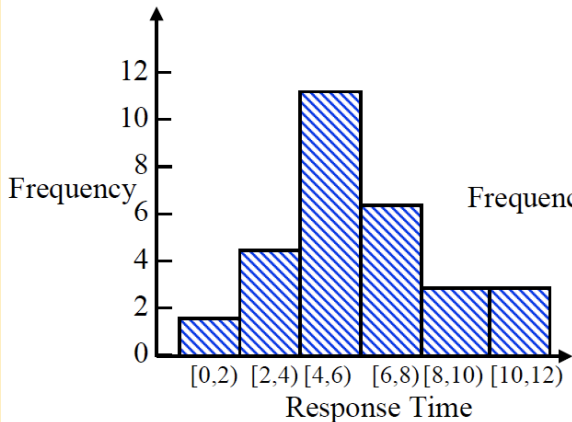
Common mistakes

Non-relevant graphic objects



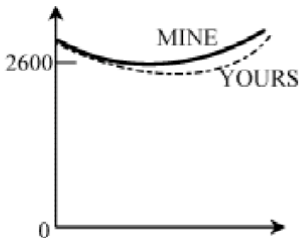
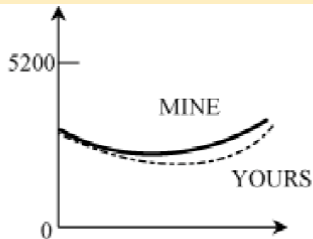
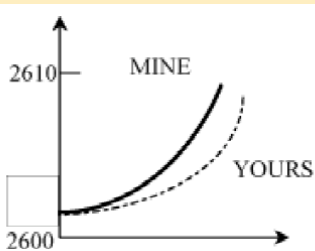
Common mistakes

Non-relevant graphic objects



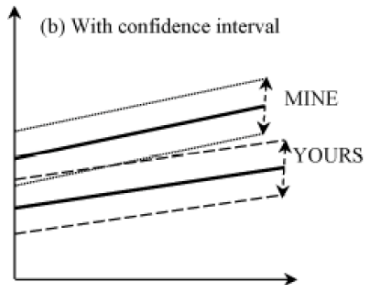
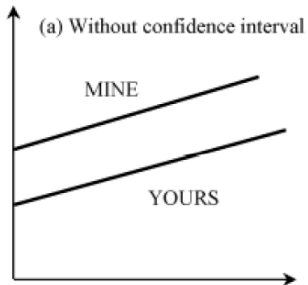
Common mistakes

Howto cheat ?



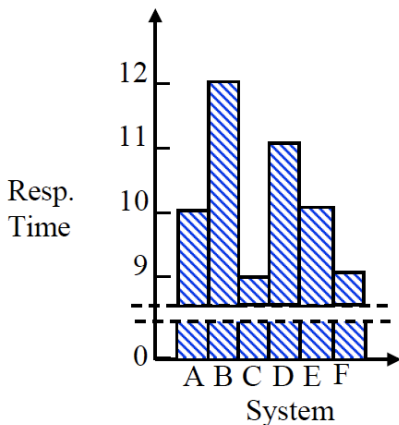
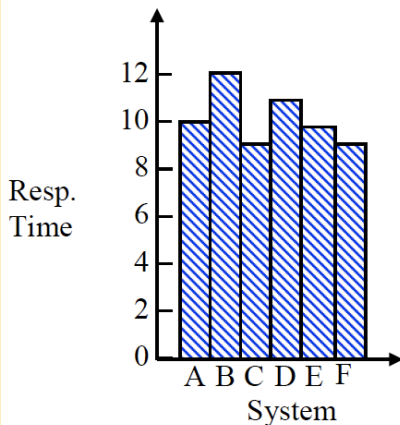
Common mistakes

Howto cheat ?



Common mistakes

Howto cheat ?



Checklist for good graphics (Jain)

- 1 Are both coordinate axes shown and labeled?
- 2 Are the axes labels self-explanatory and concise?
- 3 Are the scales and divisions shown on both axes?
- 4 Are the minimum and maximum of the ranges shown on the axes appropriate to present the maximum information.
- 5 Is the number of curves reasonably small? A line chart should have no more than six curves.
- 6 Do all graphs use the same scale? Multiple scales on the same chart are confusing. If two charts are being compared, use the same scale if possible.
- 7 Is there no curve that can be removed without reducing the information?
- 8 Are the curves on a line chart individually labeled?
- 9 Are the cells in a bar chart individually labeled?
- 10 Are all symbols on a graph accompanied by appropriate textual explanations?
- 11 If the curves cross, are the line patterns different to avoid confusion?

Checklist for good graphics (Jain)

- 12 Are the units of measurement indicated?
- 13 Is the horizontal scale increasing from left to right?
- 14 Is the vertical scale increasing from bottom to top?
- 15 Are the grid lines aiding in reading the curve?
- 16 Does this whole chart add to the information available to the reader?
- 17 Are the scales contiguous? Breaks in the scale should be avoided or clearly shown.
- 18 Is the order of bars in a bar chart systematic? Alphabetic, temporal, best-to-worst ordering is to be preferred over random placement.
- 19 If the vertical axis represents a random quantity, are confidence intervals shown?
- 20 For bar charts with unequal class interval, is the area and width representative of the frequency and interval?
- 21 Do the variables plotted on this chart give more information than other alternatives?

Checklist for good graphics (Jain)

- 22 Are there no curves, symbols, or texts on the graph that can be removed without affecting the information?
- 23 Is there a title for the whole chart?
- 24 Is the chart title self-explanatory and concise?
- 25 Does that chart clearly bring out the intended message?
- 26 Is the figure referenced and discussed in the text of the report?