



# SLAM DUNK. How to score full points on performance in SLAs.

**Adam Grummitt – [adam.grummitt@metron.co.uk](mailto:adam.grummitt@metron.co.uk)**

Director  
Metron

## Abstract: SLAs and Performance Assurance

- SLAs define IT service requirements formally
- Constrain/contract both receivers & providers
- Define/repository for Performance Targets
- Measurable key performance indicators (KPI)
- Business and application views inc. response time
- Models used to reflect and police SLAs
- Establish a performance management regime
- Threshold violations alarms and alerts
- Achieve Performance Assurance

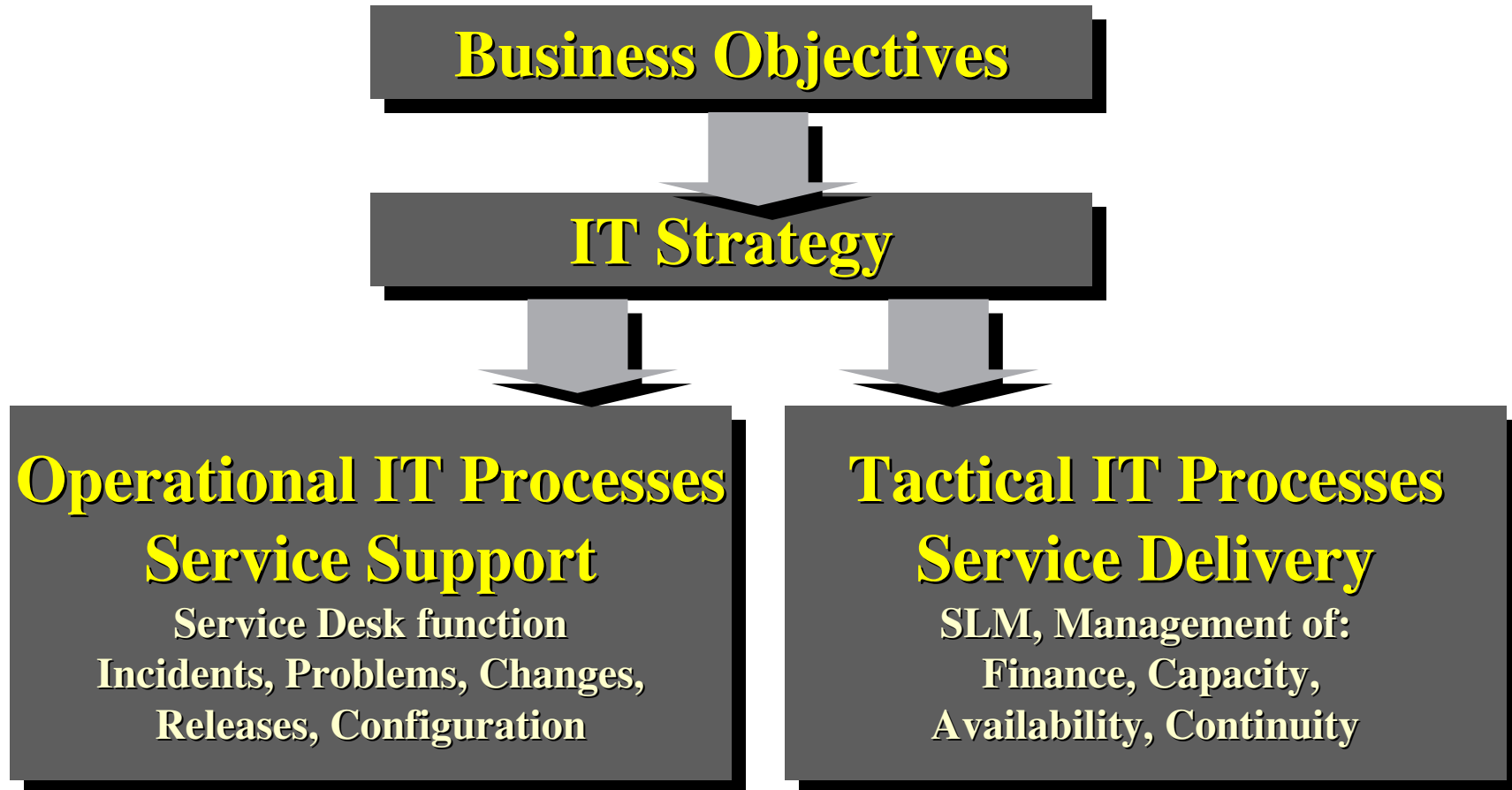
# Introduction

- SLA & SLAM – ITIL & ITSM
- Six Sigma, BSC, CoBit, QM, MOF
- Practical approach to performance in SLAs
- A skeleton SLA
- Typical outcome analysis
- Use of capacity management techniques
- Typical implementations and benefits
- Samples including ecommerce

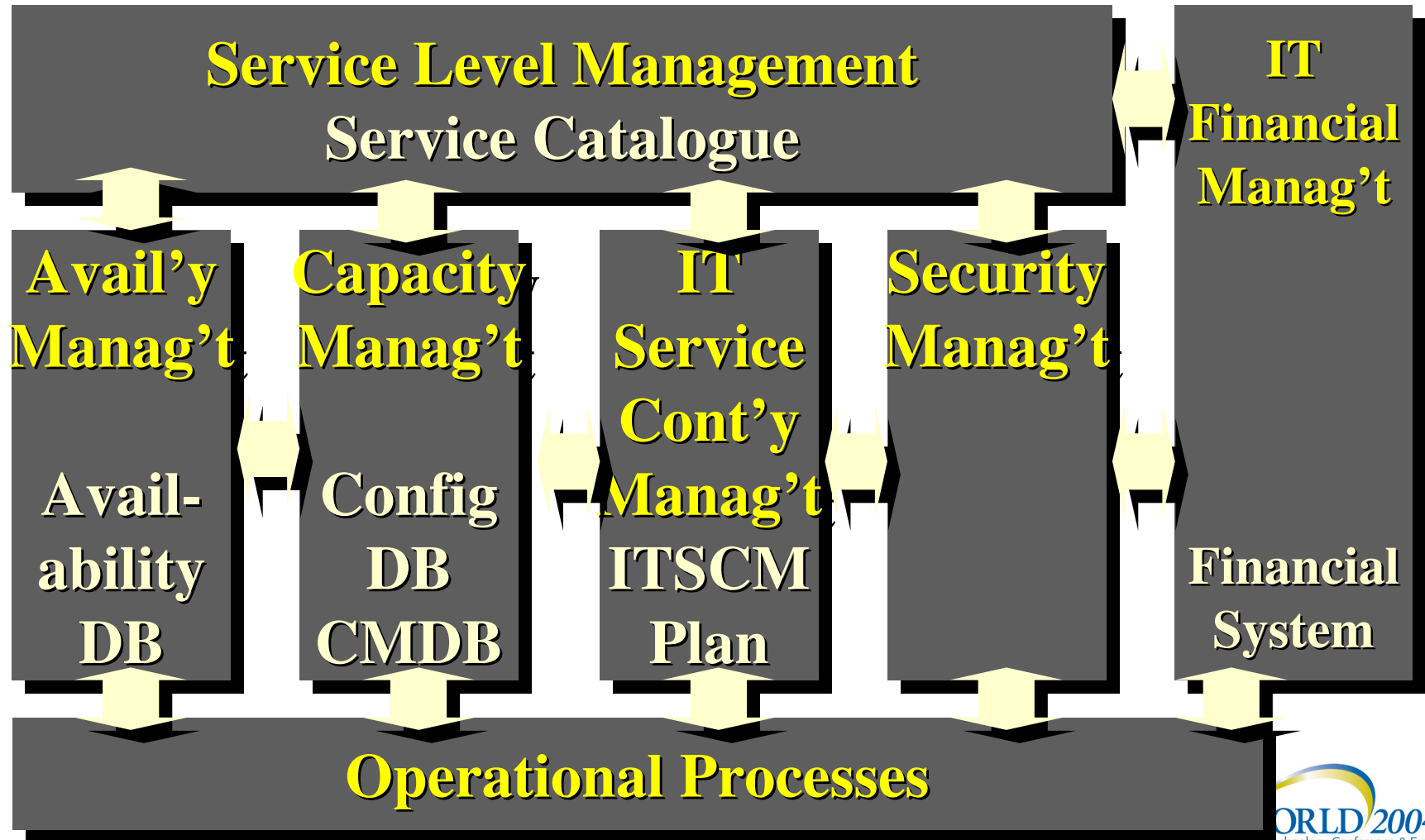
# ITIL

- The ITI Library - books & definitions
  - Service Support & Service delivery
  - Business, Infrastructure, Development, Service
- Good practice for managing IT
- Basis of BS15000, 7799 and ISO 17799 standards
- Developed by UK's OGC in the 90's
- Metron key contributor to initial Demonstrator
- itSMF
  - The IT Service Management Forum for ITIL users
  - Promotes exchange of info & experience
  - GB, NL, B, AUS, ZA, CDN, F, CH/A/D, USA

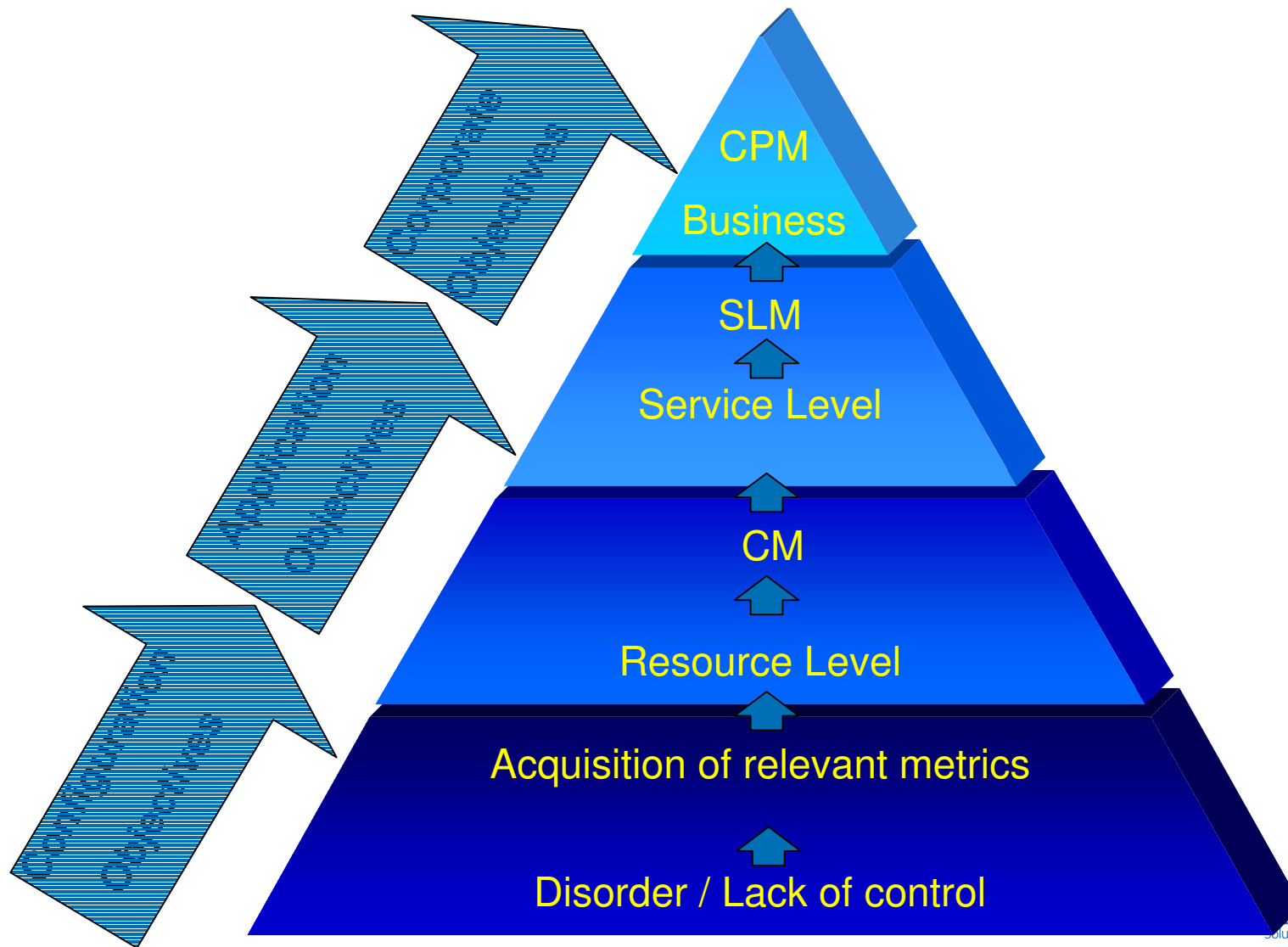
# ITIL overview



# ITIL Service Delivery Processes



# Performance Process Pyramid



# Options

## Performance Assurance

Assess improvements (EFQM)

Aim for total quality (TQM)

Operate (6 $\sigma$  - DMAIC)

People, Process (BSC)

Implement CM (ITIL)

Instal agents & create PDB

Plan attack (FMEA)

Define CPM (CMM)

Understand target (SIPOC)

Define business metrics

Performance metrics (COBIT)

Define data capture strategy

Route rationale (FMEA/Lean)

Assess position (BS15000)

Current position (ITIL/BS15000)

Understand customers (VoC)

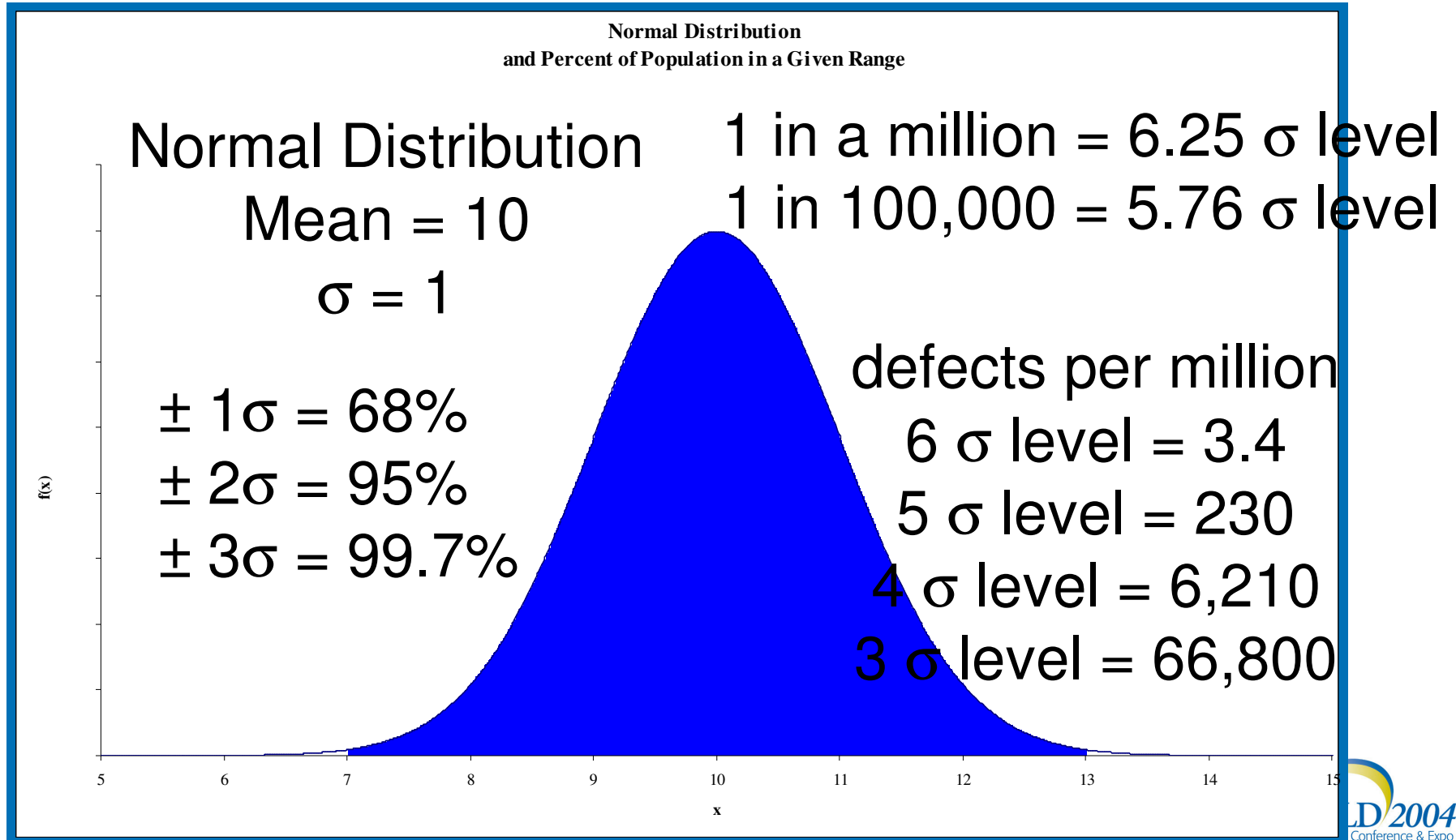
Decide on business objective

Business req's (CtQ)

Performance chaos



# Six Sigma



# BSC EG

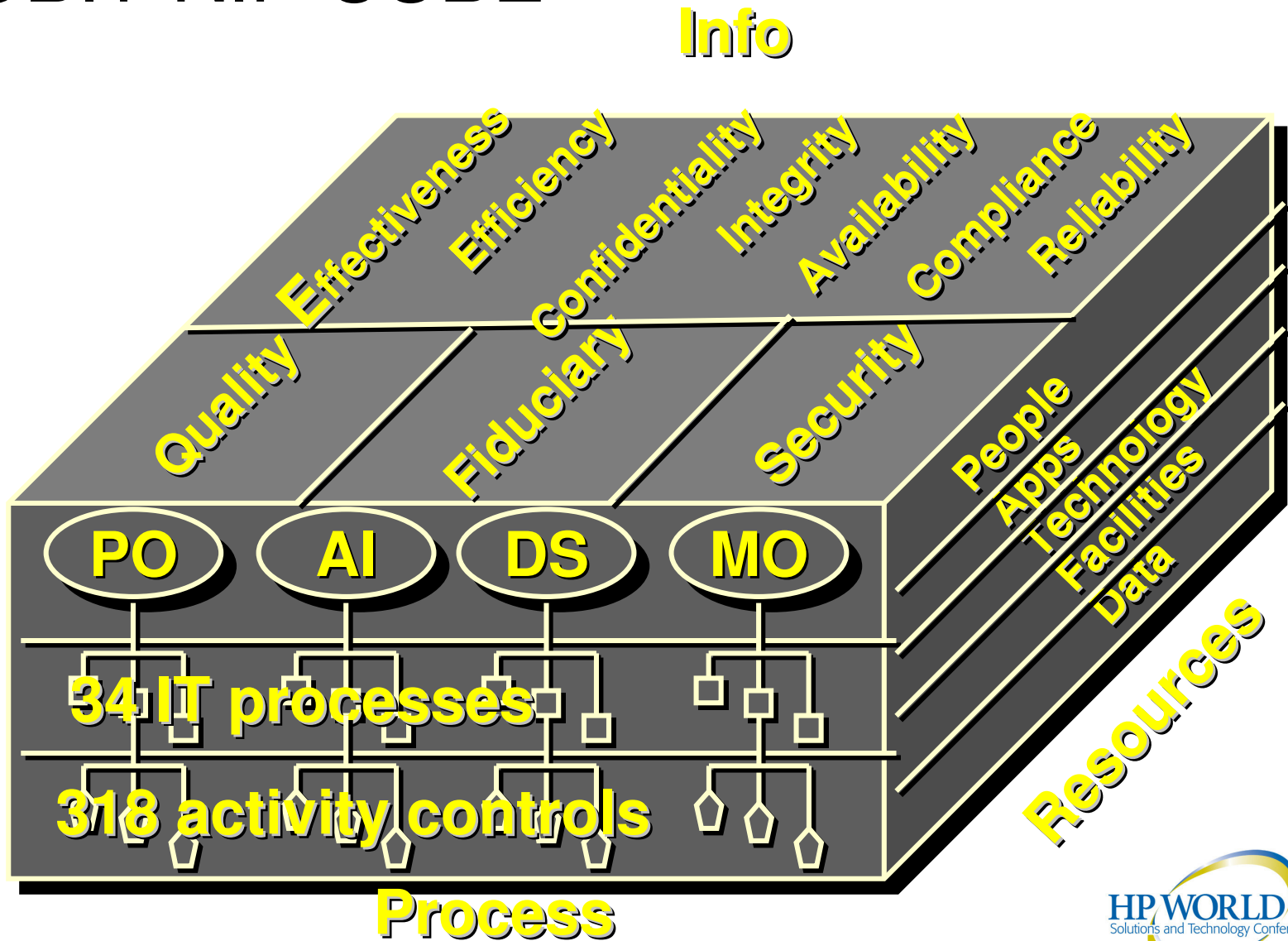
Process	KPI	#	!	$\sigma$
Finance	Cost/income	7.3		2.1
	Profit v Plan	4.2		1.8
	Risk exposure	3.2		1.6
	Earnings/share	8.8		1.2

Process	KPI	#	!	$\sigma$
Internal	Staff turnover	4.3		1.1
	Unit trans cost	9.1		0.7
	Risk exposure	4.1		1.0
	Conformance	9.4		0.3

Process	KPI	#	!	$\sigma$
Client	Intermediation	6.4		1.2
	Win rate	3.6		0.2
	Market share	4.3		0.3
	Satisfaction index	8.8		1.2

Process	KPI	#	!	$\sigma$
New	Product Sales	2.3		1.1
	Sales/employee	5.2		2.7
	Staff Attitude	3.4		1.6

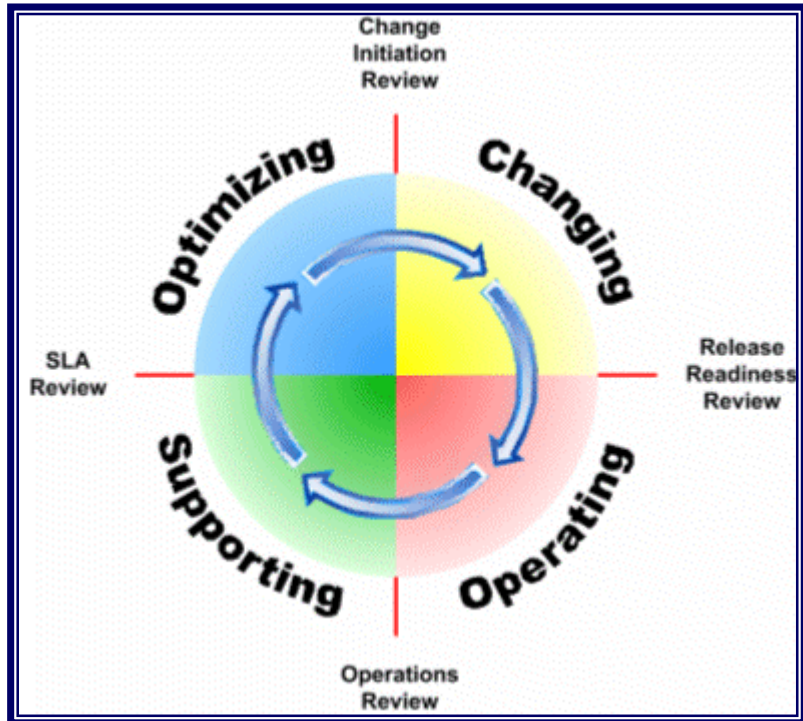
# COBIT RIP CUBE



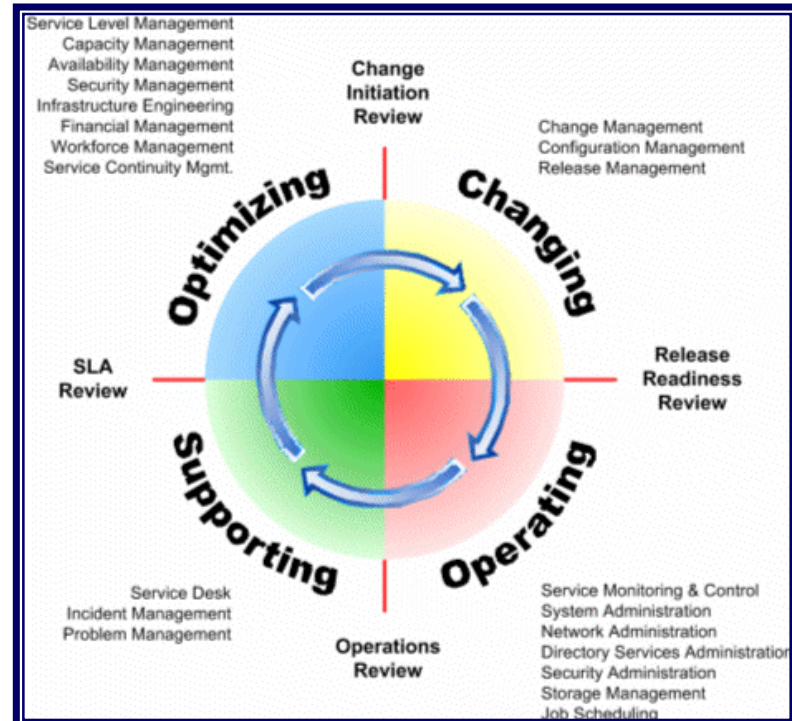
# CMM, TQM, EFQM

- CMM levels and QM
  1. Basic introduction, ad hoc solution
  2. Repeatable or reactive process
  3. Defined or proactive process
  4. Quantitatively managed process
  5. Optimised effective implementation

# MOF



Process model



Process model & ITSM functions

# SLAs

- Quantify obligations of provider & receiver
- More important if services externally charged
- Functions that the service will provide and when
- Need measurable performance indicators
- Mutual interest for it to be clear & measurable

**Contract; Repository; Yardstick**

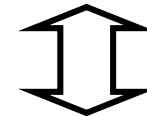
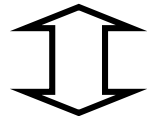


# SLAs & Capacity Management

QA ↔

Capacity Management  
(Performance Assurance)

↔ SLA



## Performance Management

- Resource accounting
- Workload balancing
- Program optimisation
- System tuning
- Alarms and alerts
- Reporting
- Tracking

## Capacity Planning

- Application sizing
- Workload trending
- Workload characterisation
- Performance Forecasting
- Modelling
- Reporting
- Tracking

# SLA processes

- Measurable numbers > arbitrary guesstimates
- Assess system at early stage in its production life
- Granularity of models  $\propto$  questions to be answered
- Split total workload into workload components
- “What-if” scenarios to assess likely bottlenecks
- Results identify thresholds for monitoring metrics
- Web reporting system - automatic alerts & alarms

**Measure: Analyse: Publish**

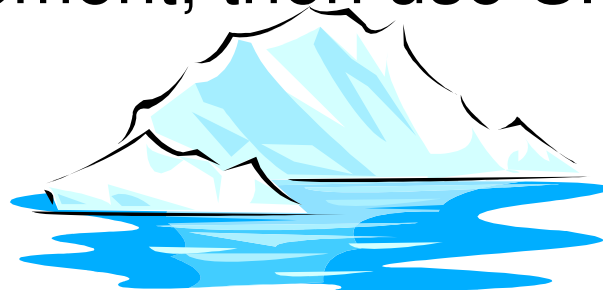


# SLA Skeleton

- Scope - parties, period, responsibilities...
- Description – application, what is (not) covered
- Service hours – normal, notice for extension...
- Service availability – % uptime in defined periods
- Service reliability – usually defined as MTBF
- User support levels – MTT respond/ resolve/ fix
- Performance – throughput, responses, turnaround
- Minimum functionality – basic service
- Contingency – continuity, security, standby
- Limitations – agreed restrictions on usage
- Financial – charging, incentives, penalties...

# SLA iceberg

- Hardware on which the system will run
- Traffic incurred
- Other workloads on the same machine
- If app on another machine/test, then measure it
- For new apps in particular, workload trials in QA
- Definition of a workload and what to measure
- Emulation or replication or a controlled workload
- If app is in development, then use SPE



# SLA & Performance

## Typical

- “Mandatory response of 3 secs; desirable 1 sec”
- “Mandatory 8 secs; desirable 5 secs for 95th %”
- “Normal maximum peak traffic of 3,600 per hour”
- “Normal service regime for 08:00 – 18:00”

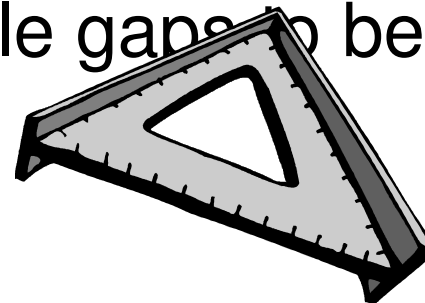
## But

- Need measures that can be monitored and used
- Spurious statistical detail re uniform distributions
- “Twice the standard deviation”, 95th percentiles

These are all part of Capacity Management

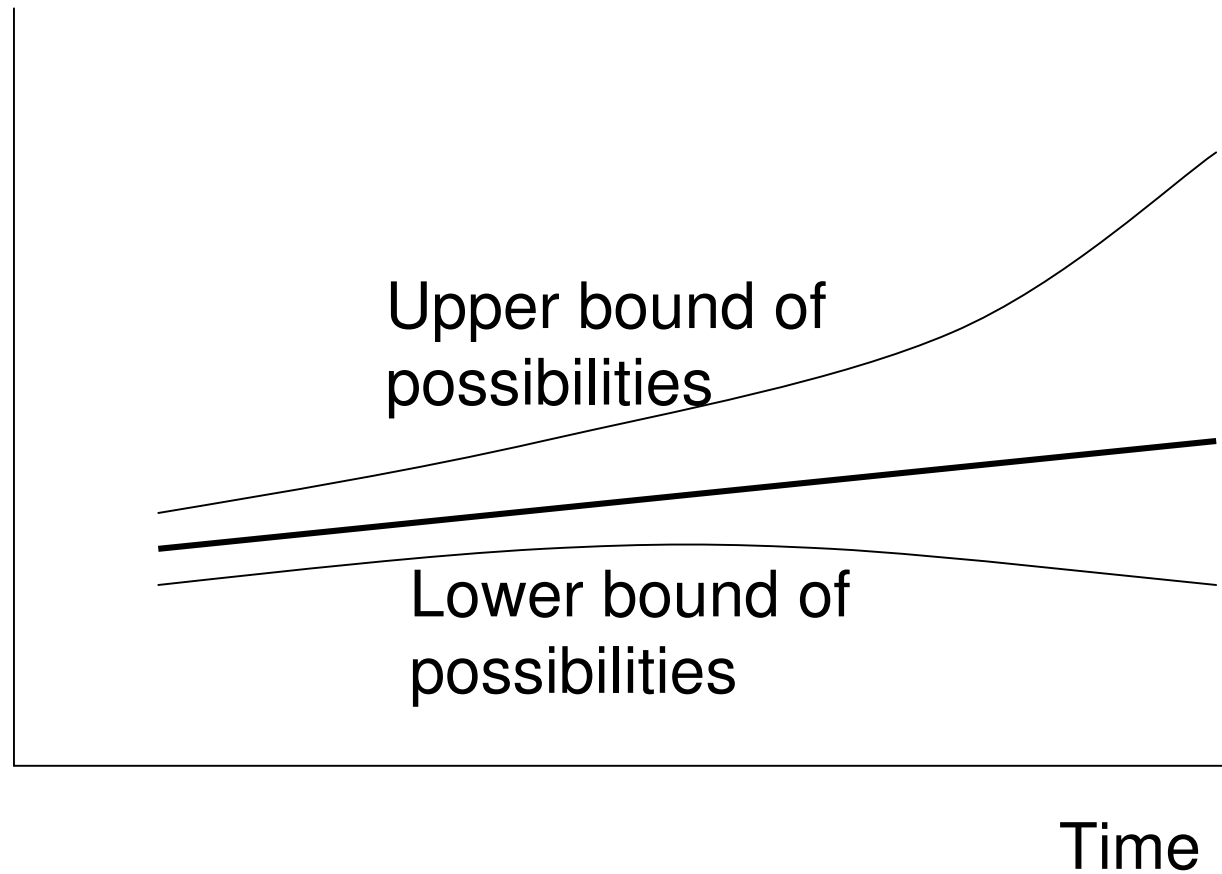
# Performance Metrics variability

- Metrics are variable in presence and reliability
- What is available is not always necessary
- What is necessary is not always available
- Both system level and user/process level
- Metrics may be sparse re mapping or responses
- Some applications are well instrumented...
- Network statistics mostly in ports, packets...
- Rules and practices enable gaps to be filled



# Sensitivity Analysis

Parameter  
Being  
forecast  
e.g.  
arrival rate,  
utilisation



# SLA outcomes

Response Time ↑

Worst

OK

Best

Agreement broken at low traffic rate

Should meet desired target at low traffic

System is performing as expected

Depends on wording of SLA

System is under pressure

May be over-configured

Prime time parameters  
Repeat for say evenings, weekends

Agreement does not apply

Depends on wording of SLA

Probably over-configured

**Mandatory**

**Desirable**

Arrival rate →

Light

Heavy

Excessive

*Normal maximum*

*Peak maximum*



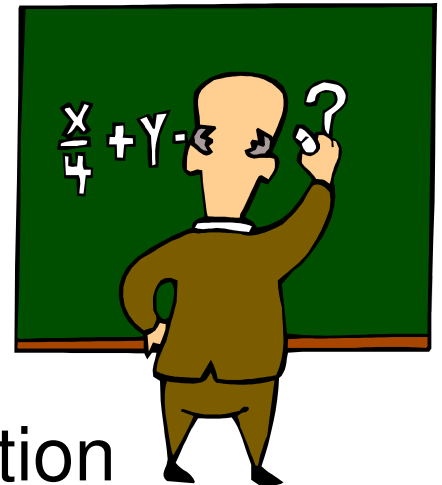
# Capacity Management & SLAM

A framework for building SLA performance:

- Characterisation of workload components
  - Evaluation of SLAs via modelling tools
  - Reporting by workload components
  - Automation of monitoring and reporting
  - Automation of alerts/alarms on violations
- Monitoring the present
  - Analysing the past
  - Predicting the future

# Analytic Model assumptions

- Use multi-class queuing network theory
- Assume large populations of transactions
- Assume exponential distributions:
  - Service times
  - Inter-arrival gaps
- “Typical” transaction is an average
- Typical SLAs assume normal distribution
- The 95th percentile usually taken as  $2\sigma$

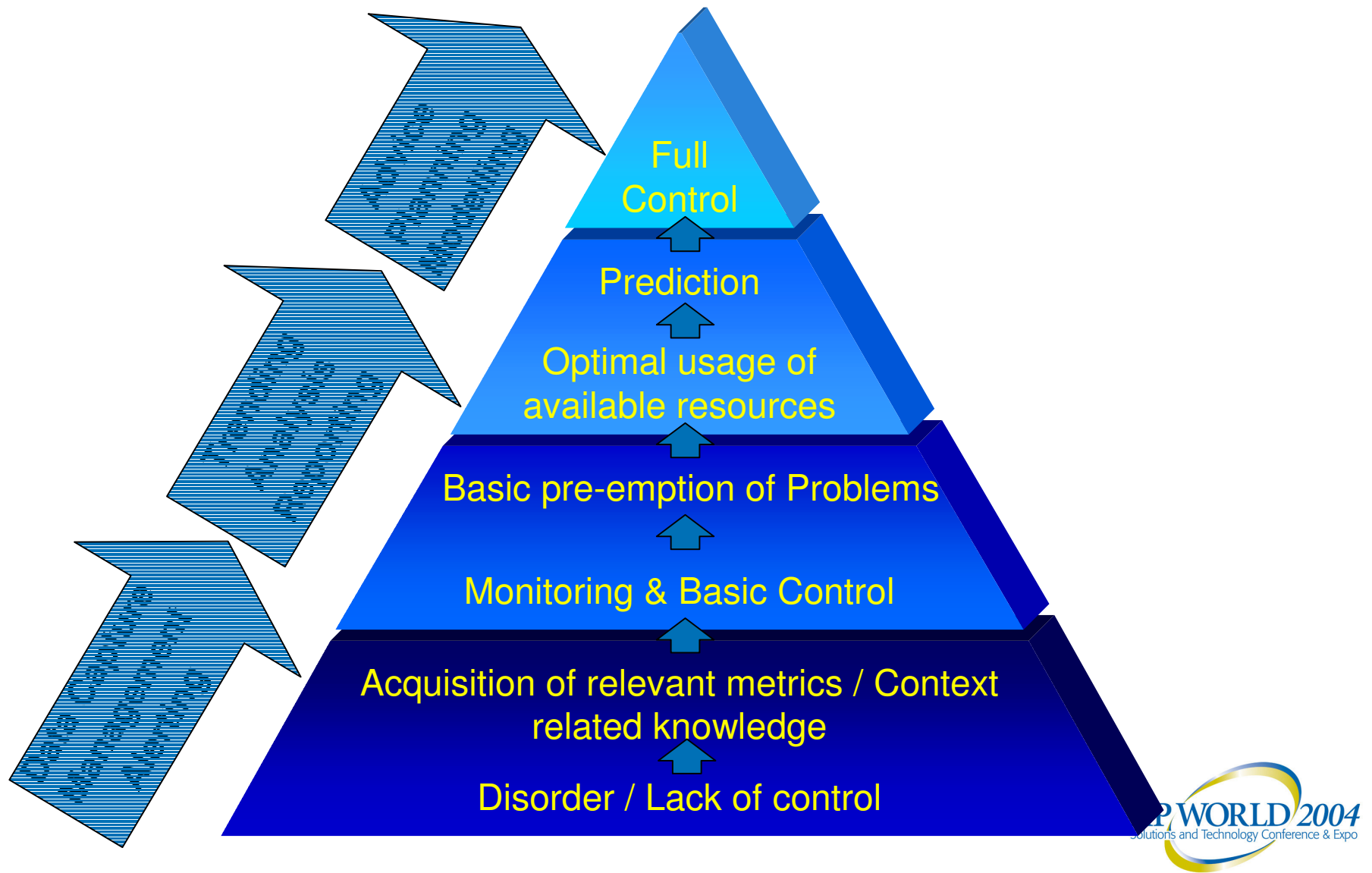




# Performance Assurance tools

- SLA definition of an app depends on the site
- Typically, n users all running a particular package
- A large number of transactions via an even larger number of processes
- Need to capture, collect and store all KPI details
- Aggregate all the resource demands for a group of processes or users = workload component
- Synthesised - usually not a “real” transaction
- Used to define a baseline situation and assess relative degradation with increasing traffic etc.

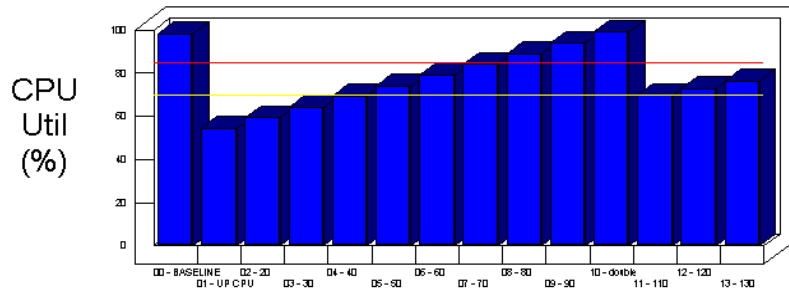
# Performance App Pyramid



CPU Utilisation Scenario Graph - ROSAGUPG



### Projected CPU Utilisation

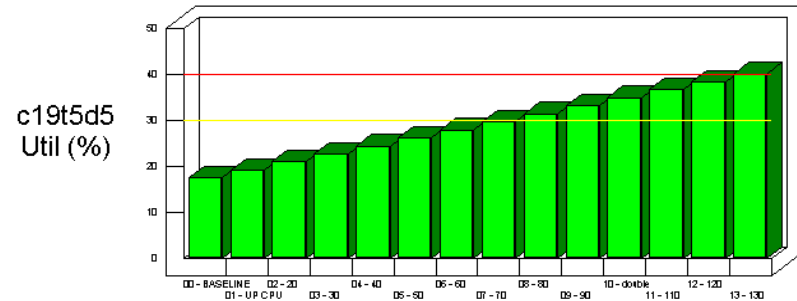


Projections

Utilisation Scenario Graph for c19t5d5 - ROSAGUPG



### Projected Device Utilisation for c19t5d5

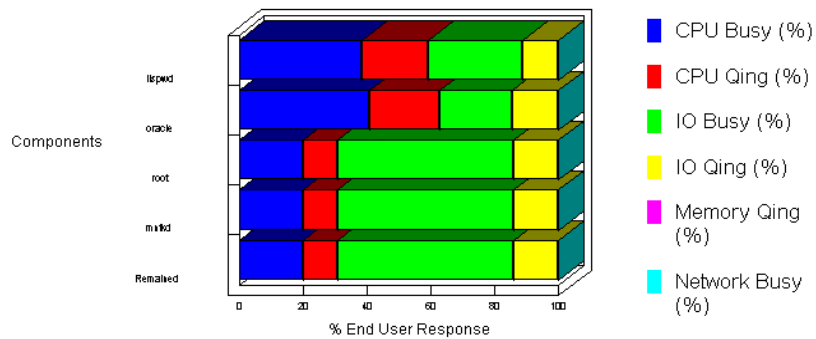


Projections

Primary Results Analysis Graph - ROSAGUPG P13



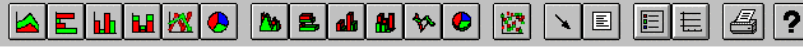
### Primary Results Analysis



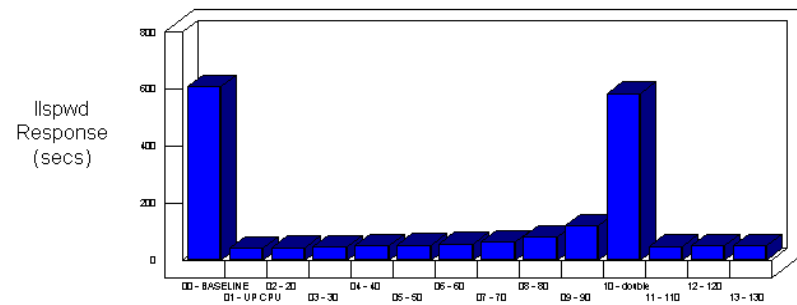
Projection: 13    Projection Name: 130

Current Model: ROSAGUPG    Current Projection: 0

llspwd Response Scenario Graph - ROSAGUPG



### Projected Response Time for llspwd



Projections



# SLAM reports

**Bulletin Name:** Head Office servers - Daily Performance  
**Installation:** Wkshop73  
**Analysis Period:** 16/05/2001 to 16/05/2001  
**Date and Time of Report:** 31/10/2001 08:54

**Status Report**

System	Oracle Database(s)	CPU	Memory	I/O	Oracle	Other
daffodil	N/A	😊	😞	😡	N/A	😞
rose	N/A	😡	😡	😡	N/A	😞
snowdrop	N/A	😊	😞	😞	N/A	😞
tulip	N/A	😊	😞	😞	N/A	😞
crocus	N/A	😊	😞	😞	N/A	😞

**Status Report**

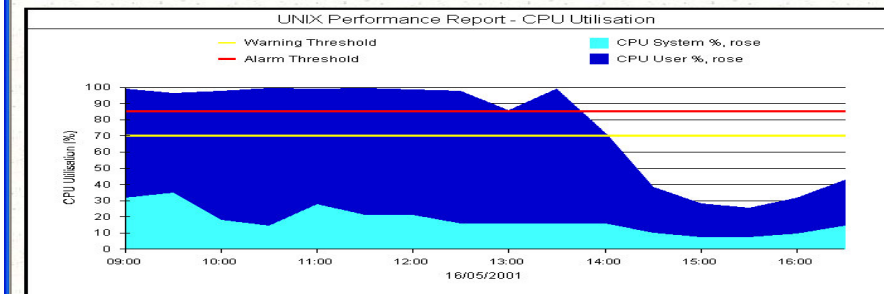
[Bulletin Index](#)

**Bulletin Name:** Head Office servers - Weekly Trend  
**Installation:** Wkshop73  
**Analysis Period:** 13/05/2001 to 19/05/2001  
**Date and Time of Report:** 31/10/2001 09:05

**Status Report**

System	Oracle Database(s)	CPU	Memory	I/O	Oracle	Other
daffodil	N/A	😞	😊	😡	N/A	N/A
rose	N/A	😊	😡	😡	N/A	N/A
snowdrop	N/A	😊	😡	😊	N/A	N/A
tulip	N/A	😊	😡	😊	N/A	N/A
crocus	N/A	😊	😊	😊	N/A	N/A

**UNIX Performance Report - CPU Utilisation**



Total CPU utilisation of system image rose

**Interpretation of UNIX Performance Report - CPU Utilisation**

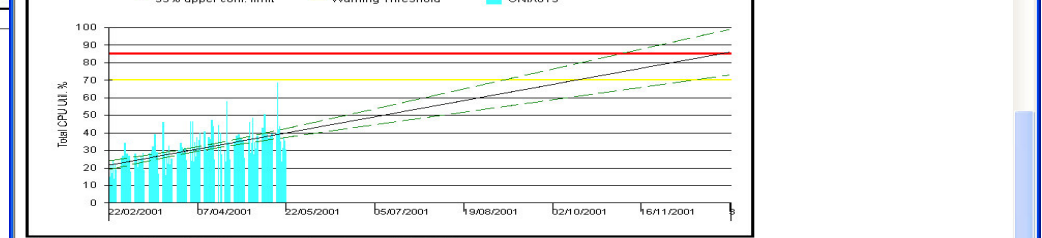
**Average value compared with thresholds**

The total CPU utilisation, averaged over the analysis period, is 75.8%. This is rather high, because it is above the lower (warning) threshold of 70.0%, but less than the (alarm) threshold of 85.0%. There is some, but not very much, spare CPU capacity to handle peaks in the workload without encountering performance problems.

**Proportion of Time compared with thresholds**

The amount of time when the total CPU utilisation is Good is 2.76 hours or 34.4% total time. The amount of time when the total CPU utilisation is Rather High is 0.5 or 6.3% of the total time. The amount of time when the total CPU utilisation is Excessively High is 4.77 hours or 59.4% of the total time.

**UNIX Trend Report - CPU Utilisation**



Trend of total CPU utilisation for system image UNIX019

**Interpretation of UNIX Trend Report - CPU Utilisation**

**Rate of Change**

The total CPU utilisation is increasing at the rate of 0.2 percent per day. This is equal to 1.4 percent per week, or approximately 6.3 percent per calendar month.

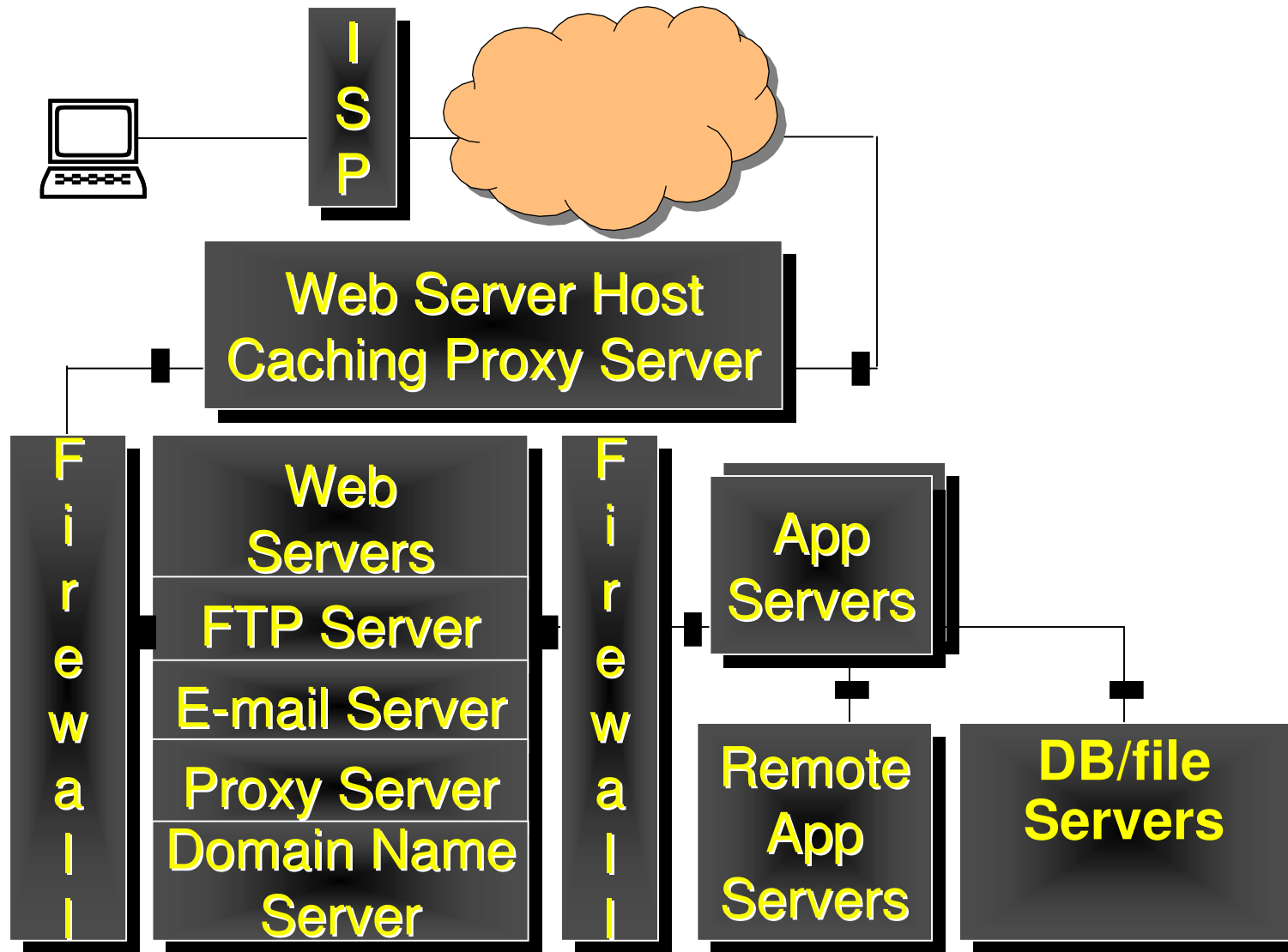
**Value at end of projection**

The projected total CPU utilisation reaches the specified value of 86.0 percent on 30/12/2001.

**Projected value compared with thresholds**

The projected total CPU utilisation rises above the upper (alarm) threshold of 85.0 percent on 25/12/2001.

# E-Commerce Multi-Tier Solution





# Multi-tier models

ecomma11 - Athene Client-Server

File View Window Help

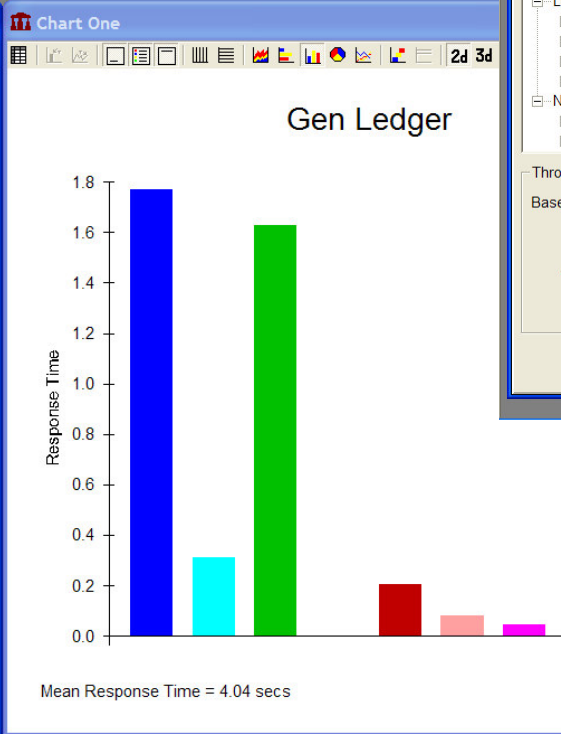
**Model Explorer**

Results Topology What If

This topology diagram shows how the nodes are connected. Right-click on the nodes to send results to the chart shown below.

1 J  
a1magbas  
Ethernet  
a1cumbas FastEthernet a1dafbas  
WAN  
a1crobas

Display on Chart - Chart One



ecomma11 - Athene Client-Server

File View Window Help

**Model Explorer**

Results Topology What If

Click on the tree below to define which model parameter you would like to alter. Use the buttons to change its value.

- a1magbas
  - Client-Server Workloads
    - GenLeqr
  - Local Workloads
    - a1magbas
    - a1cumbas
    - a1dafbas
    - a1crobas
  - Networks
    - Ethernet
    - FastEthernet

Throughput of GenLeqr  
Baseline Throughput = 1529

GenLeqr 1529

+5% Restore Edit -5% Restore All

**Chart One**

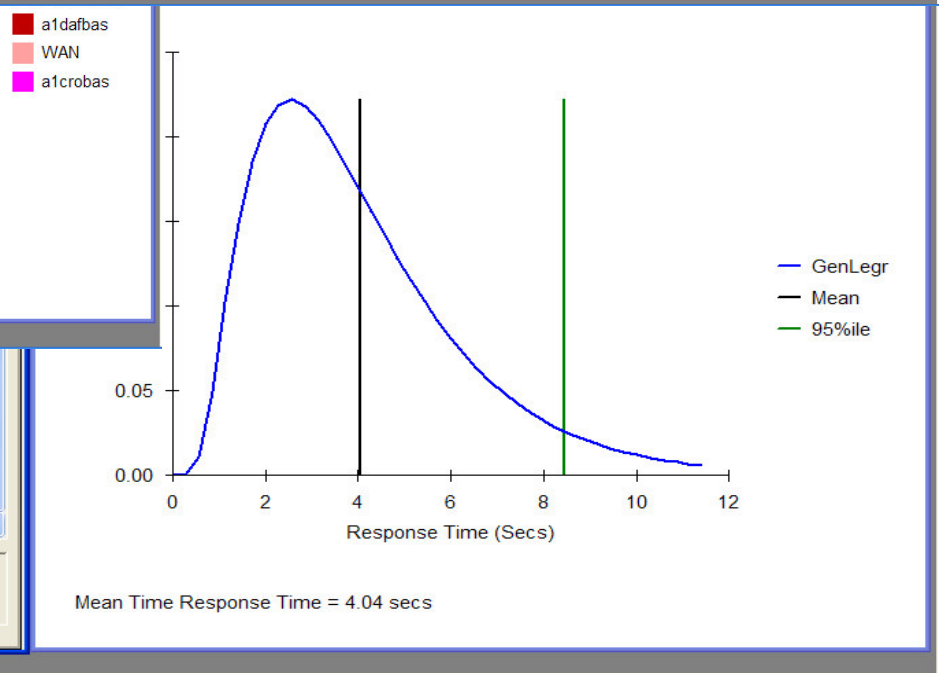
Gen Ledger steps +1% linear

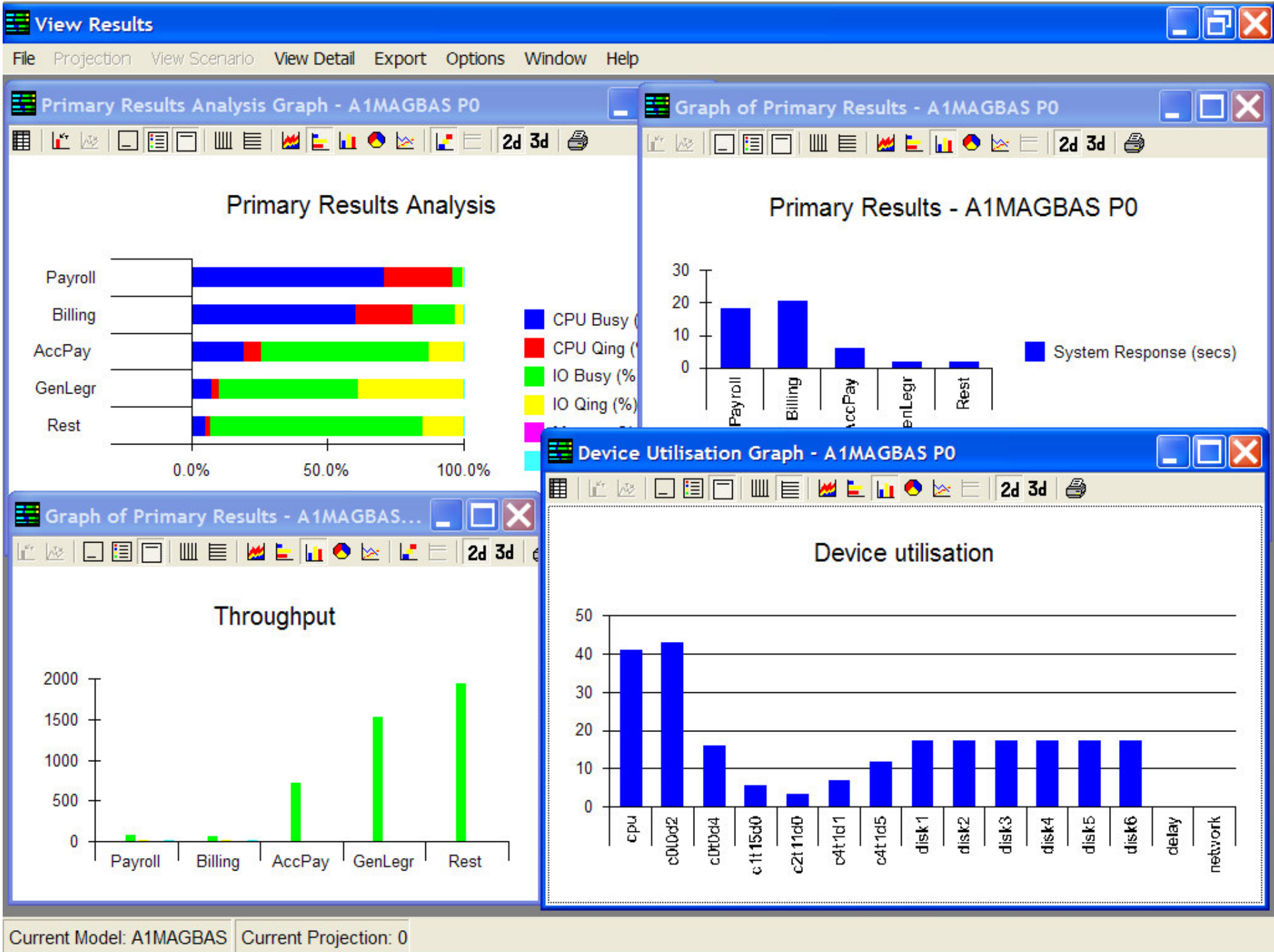
Throughput

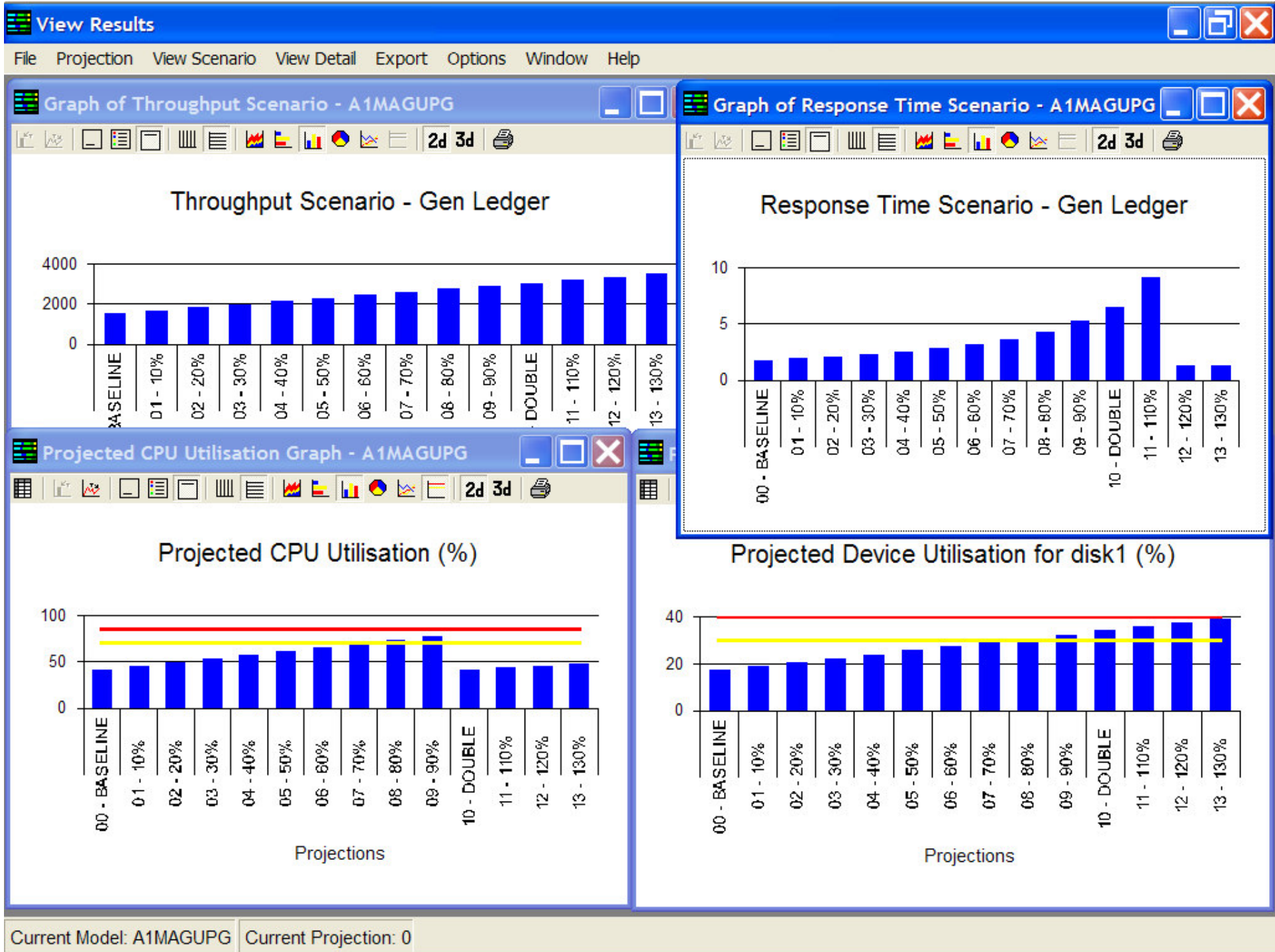
Trend Analysis

Local  
Client-Server  
GenLeqr  
Response Time Distribution  
Local  
Client-Server  
GenLeqr

Display on Chart - Chart One











**SLAM DUNK. How to score full points on performance in SLAs.**

**Adam Grummitt – [adam.grummitt@metron.co.uk](mailto:adam.grummitt@metron.co.uk)**

Director  
Metron

***Any Questions ???***

# Conclusion

- Small overhead to add performance to SLAs
- Without it, there is no performance assurance
- Only a measurable SLA can be used to police
- Modelling enables meaningful measures
- Both sides of the service have an agreed measure
- Performance of service becomes a known entity
- The service level is a sure thing; it's a SLAM dunk!



# HP WORLD 2004

Solutions and Technology Conference & Expo

Co-produced by:

