# Performance Evaluation of Computer Systems

## Course overview

Tran, Van Hoai (hoai@hcmut.edu.vn)
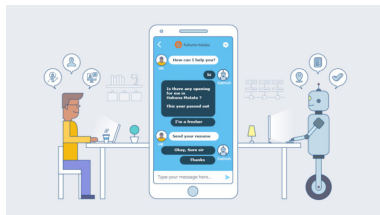
Faculty of Computer Science & Engineering
HCMC University of Technology

2023-2024/Semester 1

(source: Internet)

How to measure the quality of a commercial chatbot ?

# AI example
## Chatbot



How to measure the quality of a commercial chatbot ?

(source: Internet)

- **Self-service rate**: percentage of user sessions that did not end with a contact action after using the bot.
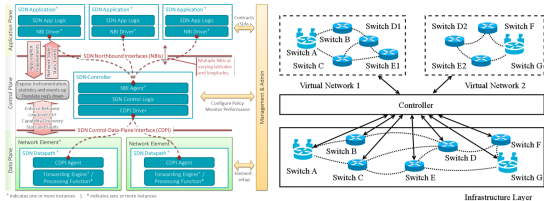- **Performance rate**: number of correct answers divided by the number of active sessions (a correct answer is an answer suggested by the bot and clicked by the user in case of multiple choices – or opened instantly in case of strong semantic matching).
- **Usage rate per login**: volume of active user sessions on the chatbot. To balance out with the average number of sessions on your website.
- **Bounce rate**: volume of sessions where the chatbot was opened but not used
- **Satisfaction rate**: average grade given when evaluating the chatbot's answers (to balance out with the evaluation rate).
- **Evaluation rate**: percentage of user sessions that have given an evaluation of the chatbot's answers at least once.
- **Average chat time**: allows you to evaluate your users' interest for your chatbot.
- **Average number of interactions**: used to evaluate the Customer Effort Score on the chatbot and must be correlated to the satisfaction rate. If the latter is very low, the bot may be engaging the users in too many branches and steps to meet their needs. In this case, a resolution can be to correct the decision trees or knowledge base architecture.
- **Goal completion rate**: in case your bot contains targeted actions like CTAs, a form or some cross-selling, that is the rate of users who have reached that specific action through the chatbot.
- **Non-response rate**: the amount of times the chatbot has failed to push some content following a user question (due to lack of content or misunderstanding).

## What is SDN? (Wiki)

SDN architectures decouple network control and forwarding functions, enabling the network control to become directly programmable and the underlying infrastructure to be abstracted from applications and network services.
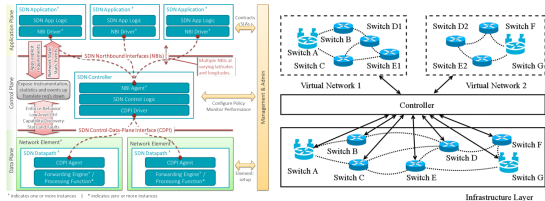
## What is SDN? (Wiki)

SDN architectures decouple network control and forwarding functions, enabling the network control to become directly programmable and the underlying infrastructure to be abstracted from applications and network services.
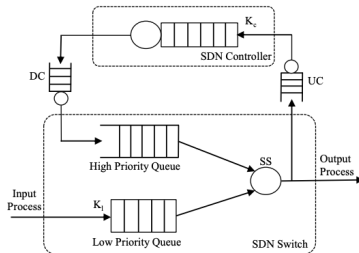
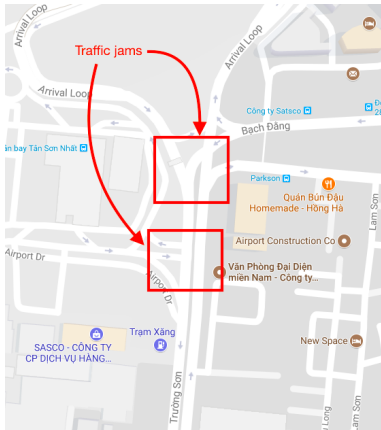Performance of SDNs can be modeled as PQ-based architecture.

### AlphaGo Zero

By playing games against itself, AlphaGo Zero surpassed the strength of AlphaGo Lee in three days by winning 100 games to 0, reached the level of AlphaGo Master in 21 days, and exceeded all the old versions in 40 days.

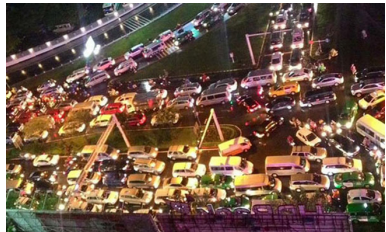- See a tracking video of isolated intersection in Vietnam
- We can train a reinforcement learning traffic light controller from scratch, and better than state-of-the-art controllers.

source: Google maps





source: cafebiz, vtv

- Investment: 242 billions VND
- Construction time: 8/2/2017 - 3/7/2017

- Investment: 242 billions VND
- Construction time: 8/2/2017 - 3/7/2017





*source: vietnamnet.vn*

- Investment: 242 billions VND
- Construction time: 8/2/2017 - 3/7/2017



source: vietnamnet.vn



source: vnexpress.vn, thanhnien.vn

Can we know which one is better
in quantitative manner
right now?

- We are here not transportation experts
  ⇒ do not discuss about which transportation methods/approaches are reasonable.
- There are something wrong in the process to choose builidng the overpass

No signalised traffic control
⇒ JAM, JAM, JAM

Signalised traffic control
⇒ ...and still JAM, JAM, JAM



### Group discussion

List possible main reasons of the jams at an intersection in the Vietnam big cities.

# Outline

## What is an engineer ?

An engineer is a professional practitioner of engineering, concerned with applying scientific knowledge, mathematics and ingenuity to develop solutions for technical problems.

(wikipedia.org)

Performance evaluation aims at predicting a system's behavior in a quantitative manner.

Performance evaluation aims at predicting a system's behavior in a quantitative manner.



source: blog.newsela.com



source: wikipedia.org

Performance evaluation aims at predicting a system's behavior in a quantitative manner.

CREATING

Putting information together in an innovative way

BLOOM'S TAXONOMY - LEARNING IN ACTION

**Examples**

- The number of stations that can be connected to a LAN and still maintain a reasonable average frame delay and throughput ?
- The fraction of calls that are blocked on outgoing lines of a company's telephone system and how much improvement we can get if an extra line is added ?
- The improvement in speedup and latency that we can achieve if we add a processor or two to a multiprocessor system ?
- The improvement in mean response time of a network if the copper wires are replaced by optical fiber ?

What is "system" in the context of performance evaluation ?

An assemblage of objects so combined by nature or human as to form an integral unit

A regularly interacting or interdependent group of objects forming a unified whole

Webster's Dictionary

A combination of components/objects that act together to perform a function not possible with any of the individual parts

IEEE Standard Dictionary of Electrical and Electronic Terms

## Two major features

1. A system consists of interacting objects/components
2. A system is associated with a function/work that it performs

Which questions are we thinking about when using the word
"**performance**"?

> Which questions are we thinking about when using the word
> "**performance**"?

- Basic concepts of work (performance metric):
    - Latency (time)
    - Bandwidth (rate)

> Which questions are we thinking about when using the word "**performance**"?

- Basic concepts of work (performance metric):
  - Latency (time)
  - Bandwidth (rate)
- How well a computer system performs a given job or activity (workload)?

> Which questions are we thinking about when using the word "**performance**"?

- Basic concepts of work (performance metric):
  - Latency (time)
  - Bandwidth (rate)
- How well a computer system performs a given job or activity (workload)?
- What is hard?
  - Performance of a computer system is multidimensional
    Complex component interaction; hard to predict how it will scale;...

CPU

arrival rate ($\lambda$) → service rate ($\mu$) →

- What is the average time it takes a job to complete service?
- What is the throughput of the system (number of jobs completed per unit time)?
- If arrival rate is doubled ($\lambda \rightarrow 2\lambda$), how much should $\mu$ increase? Do we do nothing or do we need another CPU?

arrival rate ($\lambda$)

CPU service rate ($\mu$)

- What is the average time it takes a job to complete service?
- What is the throughput of the system (number of jobs completed per unit time)?
- If arrival rate is doubled ($\lambda \to 2\lambda$), how much should $\mu$ increase? Do we do nothing or do we need another CPU?

If we need more server capacity, one queue per server what are our options?

- buy a new server with the needed capacity
- buy a few smaller servers that adds up to the required capacity
    1. one queue for all servers
    2. one queue for each server



one queue for all servers

one queue per server

# In which courses performance discussed?

- Digital systems
- Logic design with HDL
- Data structures and algorithms
- Computer architectures
- Operating systems
- Microprocessors
- Embedded systems
- Computer networks
- Distributed systems
- Management for engineers
- ...
- Graduation project

- Measurements of actual systems
- Simulations using software systems
- Mathematical modeling using techniques as queuing analysis

- **Measurements** of **actual** systems
- **Simulations** using **software** systems
- **Mathematical modeling** using techniques as **queuing analysis**

# Goals of performance evaluation

- **Compare alternative system designs**
  Example: How much memory for shortest path algorithms ?

- **Procurement**
  Example: finding a cost-effective database for a specific application ?

- **Capacity planning**
  Example: Using available resources for optimum performance

- **System tunning**
  Example: finding best set of parameters for a high performance system

- **Performance debugging**
  Example: a system runs slowly not as expected in design. Trying to detect reasons for that

- **Set expectation/Recognize relative performance**
  Example: to "guess" something not jet happen, but be useful for planning.

The course has two main goals

- Introducing a broad knowledge on performance evaluation and its application in computer science and engineering
- Presenting basic steps in performance evaluation, including
    - Performance metrics and workload
    - Analytical modeling
    - Experimental design
    - Simulation

- Specifying performance requirements
- Characterizing the load on the system (workload characterization)
- Finding the performance bottleneck (bottleneck identification)
- Determining the number and sizes of components (capacity planning)
- Evaluating design alternatives

| Part | Percentage | Assessment method |
|------|-----------|-------------------|
| Quiz 1 (Kiểm tra 1) | 20 | Multiple choice (required) |
| Quiz 2 (Kiểm tra 2) | | Multiple choice |
| Quiz 3 (Kiểm tra 3) | | Multiple choice |
| Project 1 (Bài tập lớn 1) | 30 | Analytical modeling & presentation (required) |
| Project 2 (Bài tập lớn 2) | | Python programming & presentation |
| Final exam (Thi cuối kỳ) | 50 | Multiple choice, short answer (required) |

| Part | Percentage | Assessment method |
|------|------------|-------------------|
| Quiz 1 (Kiểm tra 1) | 20 | Multiple choice (required) |
| Quiz 2 (Kiểm tra 2) | | Multiple choice |
| Quiz 3 (Kiểm tra 3) | | Multiple choice |
| Project 1 (Bài tập lớn 1) | 30 | Analytical modeling & presentation (required) |
| Project 2 (Bài tập lớn 2) | | Python programming & presentation |
| Final exam (Thi cuối kỳ) | 50 | Multiple choice, short answer (required) |

**Changes in assessment plan**

The assessment plan is tentative.

- **Study** the glossary of symbols and definitions to get to know the "language".
- **Read** ahead material before each lecture.
- **Do** lots of exercies, look at problems early from the begining of the semester.
- **Redo** your homework until it is correct.
- **To be encouraged** active learning.

- **"Art of Computer Systems Performance Analysis", R. Jain, Wiley, 1991, ISBN:0471503363** (Winner of the "1992 Best Computer Systems Book" Award from Computer Press Association")
- "Computer Systems Performance Evaluation and Prediction", Paul Fortier and Howard Michel, Digital Press, 2000.
- "Fundamentals of Performance Evaluation of Computer and Telecommunication Systems", Mohammad S. Obaidat and Noureddine A. Boudriga, Wiley, 2010.
- Internet

## Objectives

- Minimum response time (to deliver an order to the final destination.
- Average number of items delivered with this minimum response time

$\lambda_1 = \lambda^{(0)} q_1$   $\lambda_2 = \lambda^{(0)} q_2$   $q_1 + q_2 = p_1 + p_2 = p_3 + p_4 = 1$

$S_0$ indicates SOURCE

$S_1$ indicates SINK or Warehouse location (France)

L1 is in France
L2, L3 are in Morocco
L4 is in Tunisia
T indicates Transport

## Reminder

- Which algorithms have you learnt ?
- What do you remember on them in terms of characteristics ?
- How to know which one is the best ? How to compare among them ?

## Reminder

- Which algorithms have you learnt ?
- What do you remember on them in terms of characteristics ?
- How to know which one is the best ? How to compare among them ?

## Problem statement

Given a large sparse weighted graph. How to evaluate the performance of shortest path algorithms ?

- How large in number of vertices $|V|$ ?

# Shortest path problems

- How large in number of vertices $|V|$ ?
- How sparse in vertex degree ?

# Shortest path problems

- How large in number of vertices $|V|$ ?
- How sparse in vertex degree ?
- Weight function: uniform distribution in a given range $[a, b]$

# Shortest path problems

- How large in number of vertices $|V|$ ?
- How sparse in vertex degree ?
- Weight function: uniform distribution in a given range $[a, b]$
- Source vertex: full sampling (i.e., all vertices)

- How large in number of vertices $|V|$ ?
- How sparse in vertex degree ?
- Weight function: uniform distribution in a given range $[a, b]$
- Source vertex: full sampling (i.e., all vertices)

**Performance metrics**

- How large in number of vertices $|V|$ ?
- How sparse in vertex degree ?
- Weight function: uniform distribution in a given range $[a, b]$
- Source vertex: full sampling (i.e., all vertices)

**Performance metrics**

- running time

## Shortest path problems

- How large in number of vertices $|V|$ ?
- How sparse in vertex degree ?
- Weight function: uniform distribution in a given range $[a, b]$
- Source vertex: full sampling (i.e., all vertices)

**Performance metrics**

- running time
- memory consumption

# Shortest path problems

- How large in number of vertices $|V|$ ?
- How sparse in vertex degree ?
- Weight function: uniform distribution in a given range $[a, b]$
- Source vertex: full sampling (i.e., all vertices)

**Performance metrics**

- running time
- memory consumption
- Time to shortest path solution

# Shortest path problems

- How large in number of vertices $|V|$ ?
- How sparse in vertex degree ?
- Weight function: uniform distribution in a given range $[a, b]$
- Source vertex: full sampling (i.e., all vertices)

**Performance metrics**

- running time
- memory consumption
- Time to shortest path solution
- How many times to swap memory in/out (for limited memory case) ?

# Shortest path problems

- How large in number of vertices $|V|$ ?
- How sparse in vertex degree ?
- Weight function: uniform distribution in a given range $[a, b]$
- Source vertex: full sampling (i.e., all vertices)

**Performance metrics**

- running time
- memory consumption
- Time to shortest path solution
- How many times to swap memory in/out (for limited memory case) ?
- ...

# Shortest path problems

- How large in number of vertices $|V|$ ?
- How sparse in vertex degree ?
- Weight function: uniform distribution in a given range $[a, b]$
- Source vertex: full sampling (i.e., all vertices)

**Performance metrics**

- running time

## Challenge

Test performance of different one source-all destinations shortest path algorithms (i.e., Dijkstra, Bellman-Ford)
*all parameters and implementation variants should be considered*

### Discussion

What happens (for shortest path algorithms/implementations) if computer memory cannot store the data structure of the whole graph ?

## Prof. Raj Jain

Performance evaluation is an art.

## Throughputs (transactions per second) of 2 systems A and B

| System | Workload 1 | Workload 2 |
|--------|------------|------------|
| A | 20 | 10 |
| B | 10 | 20 |

## Prof. Raj Jain

Performance evaluation is an art.

## Throughputs (transactions per second) of 2 systems A and B

| System | Workload 1 | Workload 2 |
|--------|-----------|-----------|
| A | 20 | 10 |
| B | 10 | 20 |

There are 3 possible ways to compare (ratio game).

| System | Average | Average(*/B) | Average(*/A) |
|--------|---------|--------------|--------------|
| A | 15 | (2+0.5/2)=1.25 | (1+1)/2=1 |
| B | 15 | (1+1)/2=1 | (0.5+2)/2=1.25 |

Enjoy your study :-)