

# The Hong Kong Polytechnic University

## Subject Description Form

<b>Subject Code</b>	COMP5585
<b>Subject Title</b>	Cloud Computing and Hyperscale AI Infrastructure
<b>Credit Value</b>	3
<b>Level</b>	5
<b>Pre-requisite/ Co-requisite/ Exclusion</b>	Knowledge of Computer Organization and Operating Systems is preferable.
<b>Objectives</b>	<p>This is a fundamental course that provides students with the foundations of cloud computing and hyperscale AI infrastructure. It covers the principles and concepts, the technical underpinnings and supporting technologies, and the best practices and applications. The objectives of this subject include:</p> <ol style="list-style-type: none"><li>1. To provide students with a broad view of the theoretical and technological aspects that have led to the evolution of cloud computing and hyperscale AI infrastructure.</li><li>2. To teach students how cloud computing and hyperscale AI infrastructure support emerging big data and AI applications that are essential to the industry, and the requirement of working with cloud computing environments &amp; hyperscale AI infrastructure and develop cloud-based applications.</li><li>3. To equip students with the underlying technologies of cloud computing, including (1) basic knowledge of cloud computing, (2) data center networking, (3) virtualization, (4) large-scale distributed computing framework.</li><li>4. To equip students with the underlying technologies of hyperscale AI infrastructure, including (1) basic knowledge of large-scale AI model training, fine-tuning, and inference, (2) networking stacks in Artificial Intelligence Data Center, and communication libraries for large-scale AI infrastructure, such as MPI, NCCL, etc, (3) distributed strategies for model training, fine-tuning, and inference. (4) emerging system performance optimization technologies for model training, fine-tuning, and inference.</li><li>5. To equip students with the knowledge and skills for the planning, design and programming of cloud systems and AI infrastructure for the real-world applications.</li></ol>

<p><b>Intended Learning Outcomes</b></p>	<p>Upon completion of the subject, students will be able to:</p> <p><i>Professional/academic knowledge and skills</i></p> <ol style="list-style-type: none"> <li>1. Understand and appreciate the technological impact of cloud computing and hyperscale AI infrastructure for future enterprises, and the technologies underpinning it.</li> <li>2. Apply systematic and principled practices to designing, implementing and deploying cloud computing and hyperscale AI infrastructure.</li> <li>3. Review and assess the opportunities, costs, and system performance optimization for how the cloud computing support emerging AI applications.</li> </ol> <p><i>Attributes for all-roundedness</i></p> <ol style="list-style-type: none"> <li>1. Systematic and incremental approach to resolving practical enterprise computing problems and challenges.</li> <li>2. Learn to work effectively as a team member.</li> <li>3. Write technical reports and present solutions.</li> </ol>
<p><b>Subject Synopsis/ Indicative Syllabus</b></p>	<p><b>Topics:</b></p> <ol style="list-style-type: none"> <li>1. <b>Overview</b> The evolution of cloud computing paradigms; Motivations and benefits of cloud computing; Definitions and principles of cloud computing; Applications of cloud computing; The overview of hyperscale AI infrastructure.</li> <li>2. <b>Part I: Cloud Architecture and Service Models</b> Cloud architecture and major components, Physical infrastructure; Service models; Service provisioning; Representative providers and platforms (Amazon, Microsoft, IBM Google, Alibaba, etc); AWS (EC2, S3, etc.)</li> <li>3. <b>Part I: Virtualization Techniques</b> CPU virtualization; Memory virtualization; I/O virtualization; Network virtualization; Xen; Docker (Container); Kubernetes.</li> <li>4. <b>Part I: Large-Scale Distributed Computing Platform</b> Basic knowledge of large-scale distributed computing platforms, such as Hadoop, Spark, Ray, etc;</li> <li>5. <b>Part II: Networking Stacks and Communication libraries in Artificial Intelligence Data Center</b> Networking architecture in emerging AI Data Center; Communication libraries like MPI and NCCL, etc.</li> <li>6. <b>Part II: Large-Scale Model Training and Fine-Tuning</b> Basic knowledge about distributed large-scale model training and fine-tuning; Various distributed training strategies; Various model training frameworks (Pytorch, Tensorflow, etc).</li> <li>7. <b>Part II: Large-Scale Model Inference and Serving.</b></li> </ol>

	<p>Basic knowledge about distributed large-scale model inference and serving; Various distributed inference strategies; Various model inference frameworks, such as VLLM, etc.</p> <p><b>8. Part II: Emerging Optimization Technologies for Model Training, Fine-Tuning, and Inference</b></p> <p>Model Compression; Data management; Strategy design, etc.</p>																																						
<p><b>Teaching/Learning Methodology</b></p>	<p>The course is comprised of lectures, tutorials and laboratory exercises. During lecture, students are taught the important concepts and principles that drive the development of cloud computing and hyperscale AI infrastructure. In the lecture, students are encouraged to actively participate in mini-discussions and questions that are designed to reinforce their understanding of concepts taught.</p> <p>During tutorials, students will be presented with real and practical scenarios of enterprise case studies. In particular, they will be given the unique opportunities to study, analyze and propose solutions that leverage cloud computing and hyperscale AI infrastructure concepts. Small group discussions will be encouraged, and students will need to present their results and solutions in the form of reports and presentations.</p> <p>To reinforce practical aspects of their training, simple lab exercises will be conducted to expose students to the state-of-the-art tools and development environment that uses cloud computing and AI infrastructure as the underlying architecture to provide enterprise solutions.</p>																																						
<p><b>Assessment Methods in Alignment with Intended Learning Outcomes</b></p>	<table border="1" data-bbox="536 1249 1393 1749"> <thead> <tr> <th rowspan="2">Specific assessment methods/tasks</th> <th rowspan="2">% weighting</th> <th colspan="6">Intended subject learning outcomes to be assessed (Please tick as appropriate)</th> </tr> <tr> <th>a</th> <th>b</th> <th>c</th> <th>d</th> <th>e</th> <th>f</th> </tr> </thead> <tbody> <tr> <td>1. Continuous Assessments</td> <td>30%</td> <td>✓</td> <td>✓</td> <td>✓</td> <td>✓</td> <td>✓</td> <td>✓</td> </tr> <tr> <td>2. Final Examination</td> <td>70%</td> <td>✓</td> <td>✓</td> <td>✓</td> <td>✓</td> <td></td> <td></td> </tr> <tr> <td>Total</td> <td>100 %</td> <td></td> <td></td> <td></td> <td></td> <td></td> <td></td> </tr> </tbody> </table> <p><b>Explanation of the appropriateness of the assessment methods in assessing the intended learning outcomes:</b></p> <p>Students taking the subject will be assessed by performance in two parts: continuous assessments and examination. Continuous assessment may include in-class discussions / quizzes, assignments, and tests and / or exams.</p>	Specific assessment methods/tasks	% weighting	Intended subject learning outcomes to be assessed (Please tick as appropriate)						a	b	c	d	e	f	1. Continuous Assessments	30%	✓	✓	✓	✓	✓	✓	2. Final Examination	70%	✓	✓	✓	✓			Total	100 %						
Specific assessment methods/tasks	% weighting			Intended subject learning outcomes to be assessed (Please tick as appropriate)																																			
		a	b	c	d	e	f																																
1. Continuous Assessments	30%	✓	✓	✓	✓	✓	✓																																
2. Final Examination	70%	✓	✓	✓	✓																																		
Total	100 %																																						

	<p>The in-class discussions and quizzes engage students to actively participate in learning during lectures and tutorials. Students are to collaboratively work together to apply what they have learned in the class to solve practical problems. Assignment may include group projects that are designed to help students to work together in a small group to solve practical case studies and examples by applying concepts that are taught in the class. The results are to be presented in the form of reports and / or presentations. Tests and assignments are designed to help students reinforced their understanding of concepts and principles that are taught in the class. They are conducted to assess independent problem solving and critical thinking skills.</p>	
<p><b>Student Study Effort Expected</b></p>	<p><b>Class contact:</b></p>	
	<ul style="list-style-type: none"> <li>▪ Lectures, Tutorials / Labs</li> </ul>	<p>39 Hrs.</p>
	<p><b>Other student study effort:</b></p>	
	<ul style="list-style-type: none"> <li>▪ Assignment, Projects, Reading and Exam</li> </ul>	<p>66 Hrs.</p>
	<p><b>Total student study effort</b></p>	<p><b>105 Hrs.</b></p>
<p><b>Reading List and References</b></p>	<p><b><u>Reference Books for Cloud Computing:</u></b></p> <ol style="list-style-type: none"> <li>1. “Essentials of Cloud Computing”, by Chellammal Surianarayanan and Pethuru Raj Chelliah. Publisher: Springer, edition: 2019.</li> <li>2. “Cloud Computing Solutions Architect: A Hands-On Approach”, by Arshdeep Bahga and Vijay Madisetti. Arshdeep Bahga &amp; Vijay Madisetti, 2019.</li> <li>3. Articles from web, technical journals, and conference proceedings will be handed out or posted on L@PU Blackboard when needed.</li> </ol> <p><b><u>Reading List for AI Infrastructure:</u></b></p> <ol style="list-style-type: none"> <li><b>1. Model Training and Fine-Tuning:</b> <ol style="list-style-type: none"> <li>(1) PipeDream: Generalized Pipeline Parallelism for DNN Training (SOSP'19)</li> <li>(2) Alpa: Automating Inter- and Intra-Operator Parallelism for Distributed Deep Learning (OSDI'22)</li> </ol> </li> <li><b>2. Model Inference:</b> <ol style="list-style-type: none"> <li>(1) Orca: A Distributed Serving System for Transformer-Based Generative Models (OSDI'22)</li> <li>(2) Efficient Memory Management for Large Language Model Serving with PagedAttention (SOSP'23)</li> <li>(3) <a href="https://github.com/xlite-dev/Awesome-LLM-Inference?tab=readme-ov-file#Mixture_of_Experts_LLM_Inference">https://github.com/xlite-dev/Awesome-LLM-Inference?tab=readme-ov-file#Mixture_of_Experts_LLM_Inference</a></li> </ol> </li> </ol>	