# Cloud Computing Final Project Demo
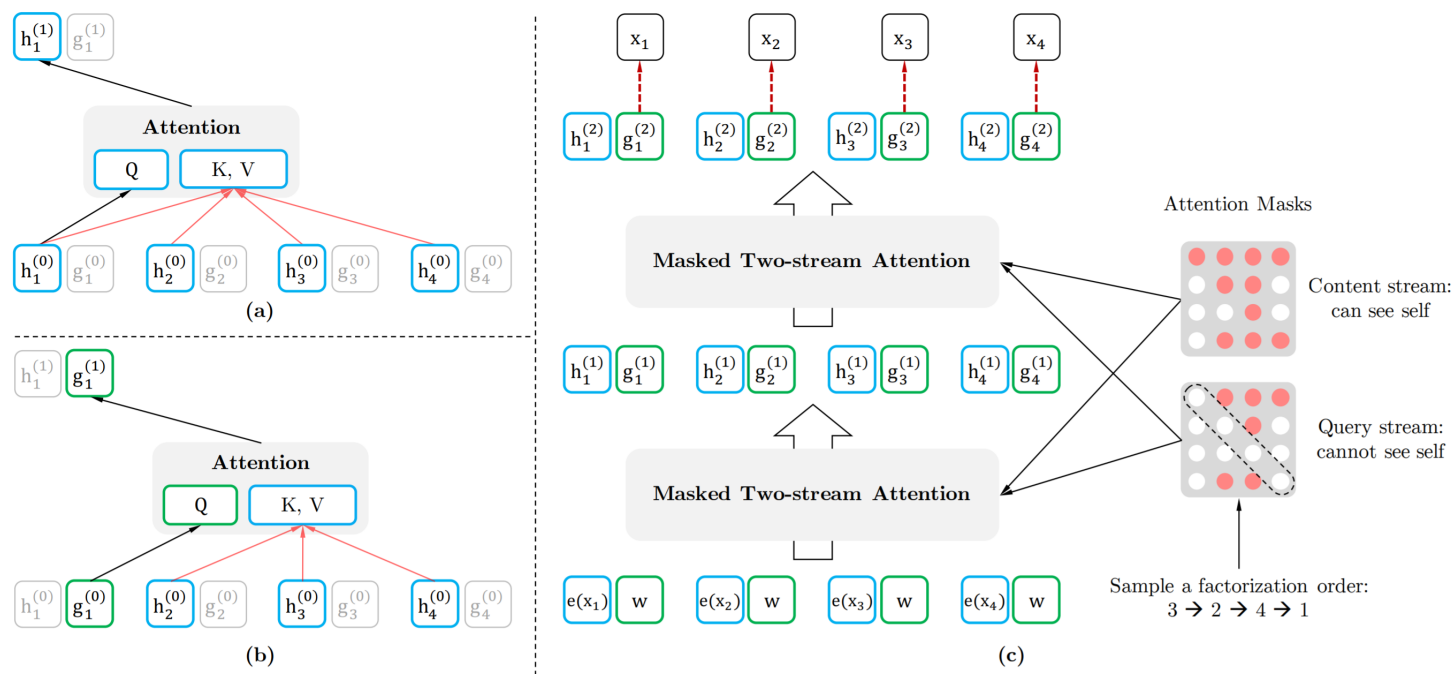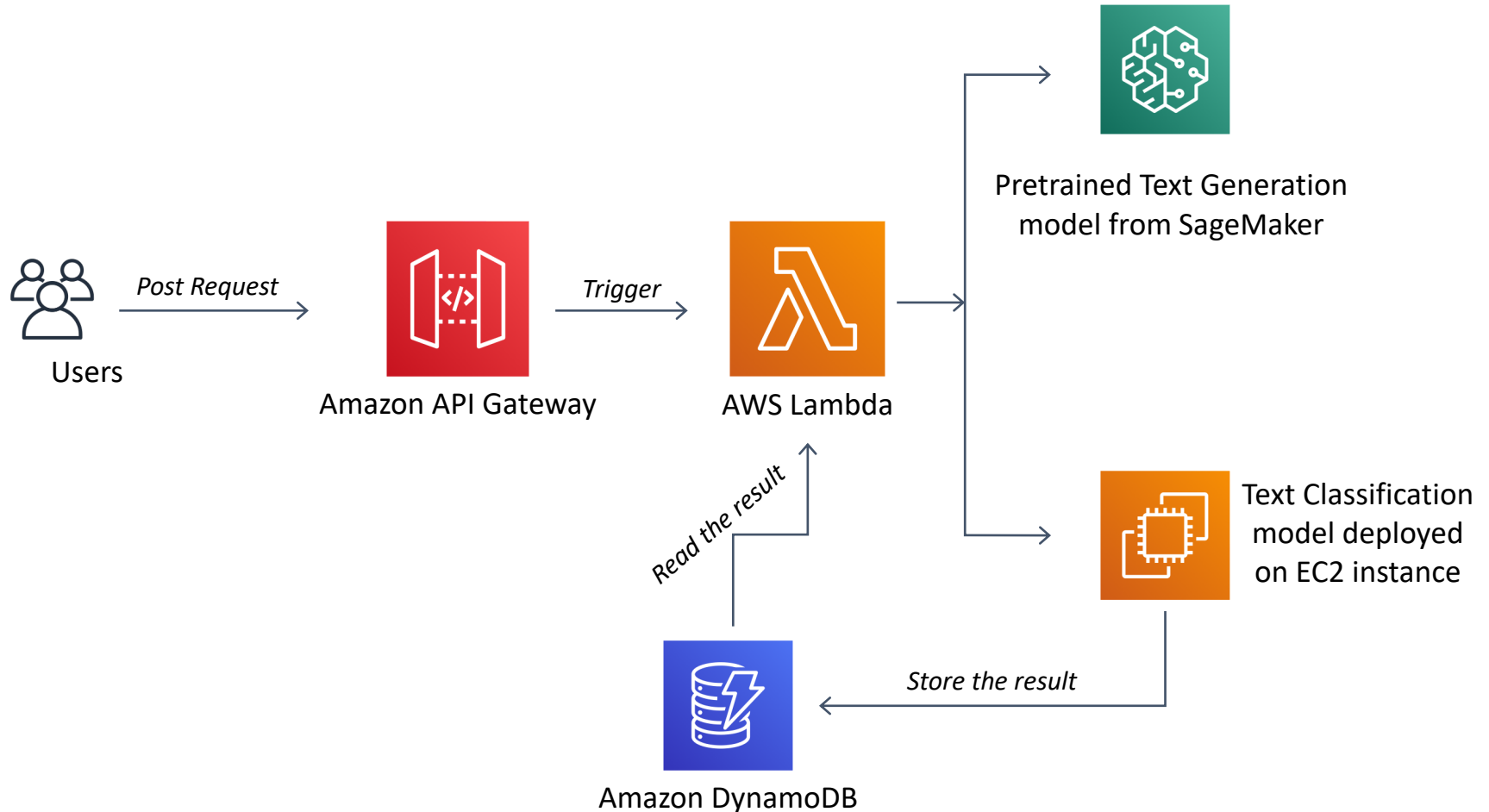
Team 08

# Model description (Fine tuned XLNet)



Figure 2: (a): Content stream attention, which is the same as the standard self-attention. (b): Query stream attention, which does not have access information about the content $x_{z_t}$. (c): Overview of the permutation language modeling training with two-stream attention.

# System Structure (Ideal)

| | id | result | |
|---|---|---|---|
| ☐ | 29146 | 0 | |
| ☐ | 162417 | 0 | |
| ☐ | 226080 | [1] | |
| ☐ | 245468 | [0] | |
| ☐ | 258456 | 1 | |
| ☐ | 497517 | [0] | |
| ☐ | 529481 | 1 | |
| ☐ | 609092 | 1 | |
| ☐ | 628218 | [0] | |
| ☐ | 750054 | 1 | |
| ☐ | 792284 | 1 | |

# gpt2-demo

## Endpoint settings

Name

**gpt2-demo**

Status

⊘ InService

# Q1: What's the purpose of the database?

- EC2 instance is not always available.
- The Lambda would first need to **turn on** the instance when requests arrive, and then sends a shell script to the instance (through SSM client).
- The script would
  - Activate the proper environment
  - Run the classification script
- The script will store the result in DynamoDB with a key from the Lambda.

# Q2: Why not perform classification inside the Lambda

- XLNet is a HUGE model, which cannot easily fit into Lambda's limited space.

- Deploying `PyTorch`, `Transformers` on Lambda has some tricky dependency problems.

- Lambda does not have GPU.

# Practical Issue

- For some reasons, we cannot create IAM role to make components interact with each other.

- So we replace the Lambda and the API Gateway with another EC2 instance, which runs a Flask server.

- Meanwhile, SSM agent is replaced by another Flask server hosted on inference server.

# Fake Lambda       =_=

```
setting up tree (1.7.0 5) ...
(base) ubuntu@ip-172-31-23-131:~/lambda$ tree
.
|-- classfication_callback.py
|-- __pycache__
|   |-- classfication_callback.cpython-37.pyc
|   `-- sagemaker_callback.cpython-37.pyc
|-- sagemaker_callback.py
`-- server.py

1 directory, 5 files
```

# Future work

- Deploy custom classification task on SageMaker

- Accelerate/Scale up Giant Language Model inference.

- Taking real-world issues into account, e.g. high concurrency.