# Experimentation Setup



DL Model

Measure Improvement

# AWS Lambda

- The function



- CLI integration with local system to run function invocation scripts

```
sid@sids-machine:~$ aws lambda invoke --function-name cOLD --cli-binary-format
raw-in-base64-out --payload '{"length": 5}' output.txt
$LATEST 200
```

- Python Script for invocation

```python
import subprocess
import time

def invoke_lambda():
    command = [
        "aws", "lambda", "invoke",
        "--function-name", "cOLD",
        "--cli-binary-format", "raw-in-base64-out",
        "--payload", '{"length": 5}',
        "output.txt"
    ]
    subprocess.run(command)

def main():
    offset = 300
    interval = 60  # Interval in seconds
    num_iterations = 5  # Number of iterations
    for i in range(num_iterations):
        invoke_lambda()
        print(f"Invocation {i+1} completed.")
        time.sleep(offset + (interval * (i + 1)))  # Increasing interval with each iteration

if __name__ == "__main__":
    main()
```

# AWS Lambda

- 1-D grid search for identifying cold -start
- Adding Cloudwatch for logs with appropriate parameters

Cold-start Indicator

| # | Timestamp | RequestId | LogStream | Duratio… | BilledD… | Initialization |
|---|-----------|-----------|-----------|----------|----------|----------------|
| ▶ 1 | 2024-04-01T20:19:37.856+05:30 | 66b67234-5606-40dc-b1… | 2024/04/01/[$LATEST]19111… | 8.76 | 9.0 | |
| ▶ 2 | 2024-04-01T20:14:16.716+05:30 | 154bb372-5c01-4622-84… | 2024/04/01/[$LATEST]19111… | 17.34 | 18.0 | 79.43 |
| ▶ 3 | 2024-04-01T20:09:04.499+05:30 | 1504d787-9066-4a0e-8c… | 2024/04/01/[$LATEST]a890f… | 1.51 | 2.0 | |
| ▶ 4 | 2024-04-01T20:04:03.436+05:30 | 7b6b9dd1-2b77-453b-95… | 2024/04/01/[$LATEST]a890f… | 4.23 | 5.0 | |
| ▶ 5 | 2024-04-01T19:59:12.428+05:30 | 7f5f589c-b84d-4984-98… | 2024/04/01/[$LATEST]a890f… | 1.76 | 2.0 | |
| ▶ 6 | 2024-04-01T19:59:05.279+05:30 | 6142fe4c-b7f9-4a23-8e… | 2024/04/01/[$LATEST]a890f… | 12.69 | 13.0 | |
| ▶ 7 | 2024-04-01T19:53:55.668+05:30 | a93ebeda-02ff-4343-b0… | 2024/04/01/[$LATEST]a890f… | 2.27 | 3.0 | 83.44 |
| ▶ 8 | 2024-04-01T19:47:53.601+05:30 | fc415174-3cb5-4283-8f… | 2024/04/01/[$LATEST]3f686… | 7.88 | 8.0 | |
| ▶ 9 | 2024-04-01T19:47:32.688+05:30 | a7bece85-3875-48ae-98… | 2024/04/01/[$LATEST]3f686… | 2.01 | 3.0 | 84.39 |

Container teardown time is ~ 300 seconds

# Log-monitoring

# Results

- For periodic Invocation (Time period = 6 minutes)
  - Test-time (3+ hours)



Comparison for periodic invocation

# Results

- For periodic Invocation
  - Test-time (3+ hours)

|  | Without Prewarming(in ms) | With Pre-warming(in ms) | improvement | Extra billing per invocation(average) | Mishits per invocation |
|---|---|---|---|---|---|
| Periodic (Average execution time) | 95.23 | 6.21 | 15x | | |
| Periodic (Median execution time) | 90.83 | 2.1 | 45x | 6.25 ms | 0 |

# Results

● For sinusoidal invocation with thresholding
  ○ Test-time (3+ hours)



Comparison for Sinusoidal

# Results

- For sinusoidal invocation with thresholding
  - Test-time (3+ hours)

|  | Without Prewarming(in ms) | With Pre-warming(in ms) | improvement | Extra billing per invocation(average) | Mishits per invocation |
|---|---|---|---|---|---|
| Periodic (Average execution time) | 65.59 | 8 | 8x | | |
| Periodic (Median execution time) | 85.18 | 4.32 | 20x | 6.36 ms | 0.03 (6 in 180) |

# Results

- For wikipedia traffic data with thresholding
  - Test-time (3+ hours)


comparison_plot_wiki_views

# Results

- For wikipedia traffic data with thresholding
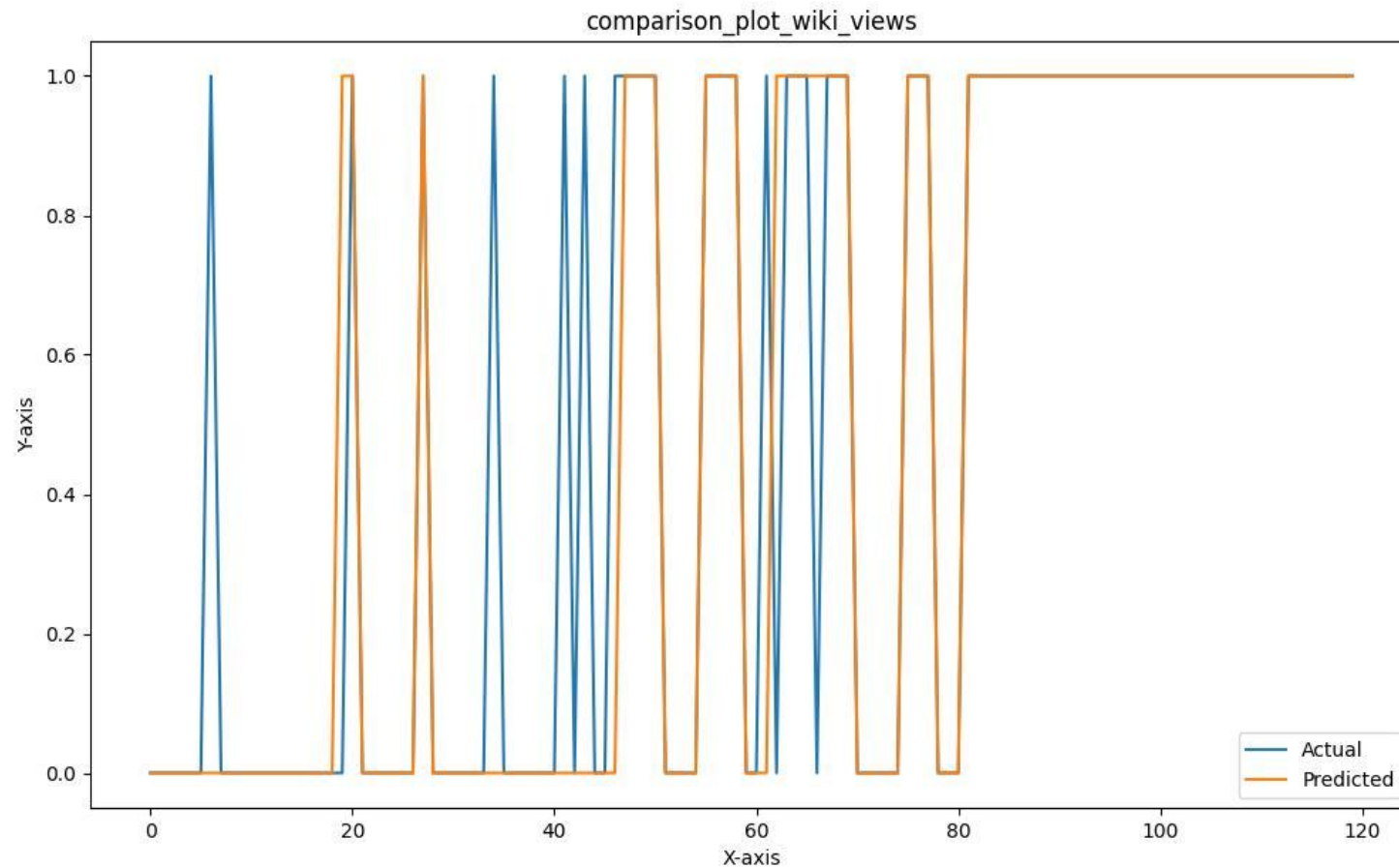  - Test-time (2+ hours)

|  | Without Prewarming(in ms) | With Pre-warming(in ms) | improvement | Extra billing per invocation(average) | Mishits per invocation |
|---|---|---|---|---|---|
| Periodic (Average execution time) | 12.12 | 6.69 | 2x | | |
| Periodic (Median execution time) | 1.55 | 1.56 | - | 3.07 ms | 0.06 |

# Conclusion & Future Scope

- Server prewarming via time-series forecasting gives a significant improvement in reducing the function execution time.
- Especially useful where the micro-service architecture uses function invocations in series i.e the output of one is used by the next until the user request is met
- A drawback of this solution is the extra billing per invocation which leads to higher costs

Future Scope
- Server - aware prewarming to reduce extra billing per invocation i.e even if the program anticipates an invocation it should send the prewarming request only if the function container is in cold state
  - Requires server side integration

# Thank you !