



# Go & Python Event driven ML

Lalafo case

Новые функции

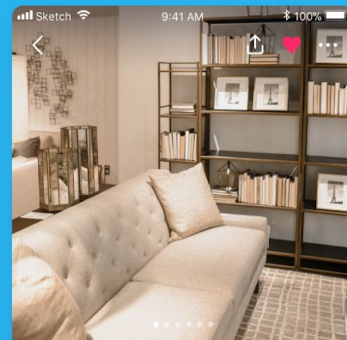
# Искусственный интеллект помогает продавать

lalafo использует собственную технологию компьютерного зрения и распознает товары по фотографии.



Bicycle

Loudspeakers



82 000 USD

## Cosy apartment in heart of Baki

A cosy and typical apartment in the heart of Baki. It is newly renovated by an interior designer. Very comfortable... [show more](#)

Furnished: Yes

Bedrooms: 2 bedroom

Area: 80

45 ❤️ 12

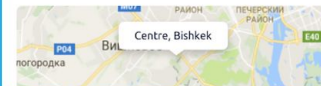
25 Sep 2017



**Auto Salon** PRO

Online

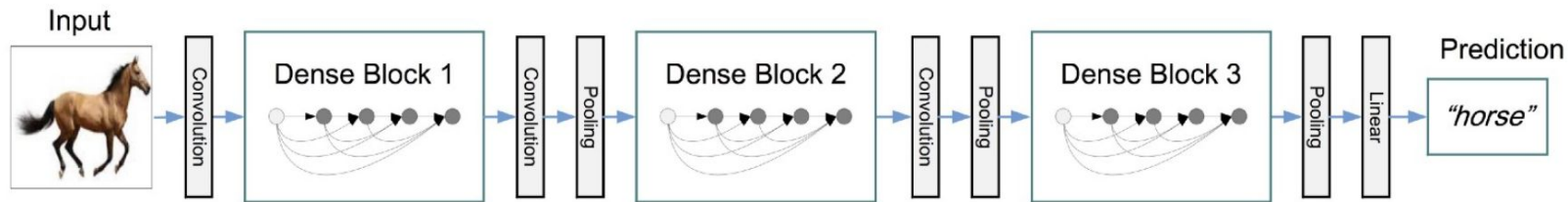
The average response speed 15 minutes



# How it all started



Part 1. The Phantom Menace



# PYTORCH



python™



Tornado



amazon  
web services™

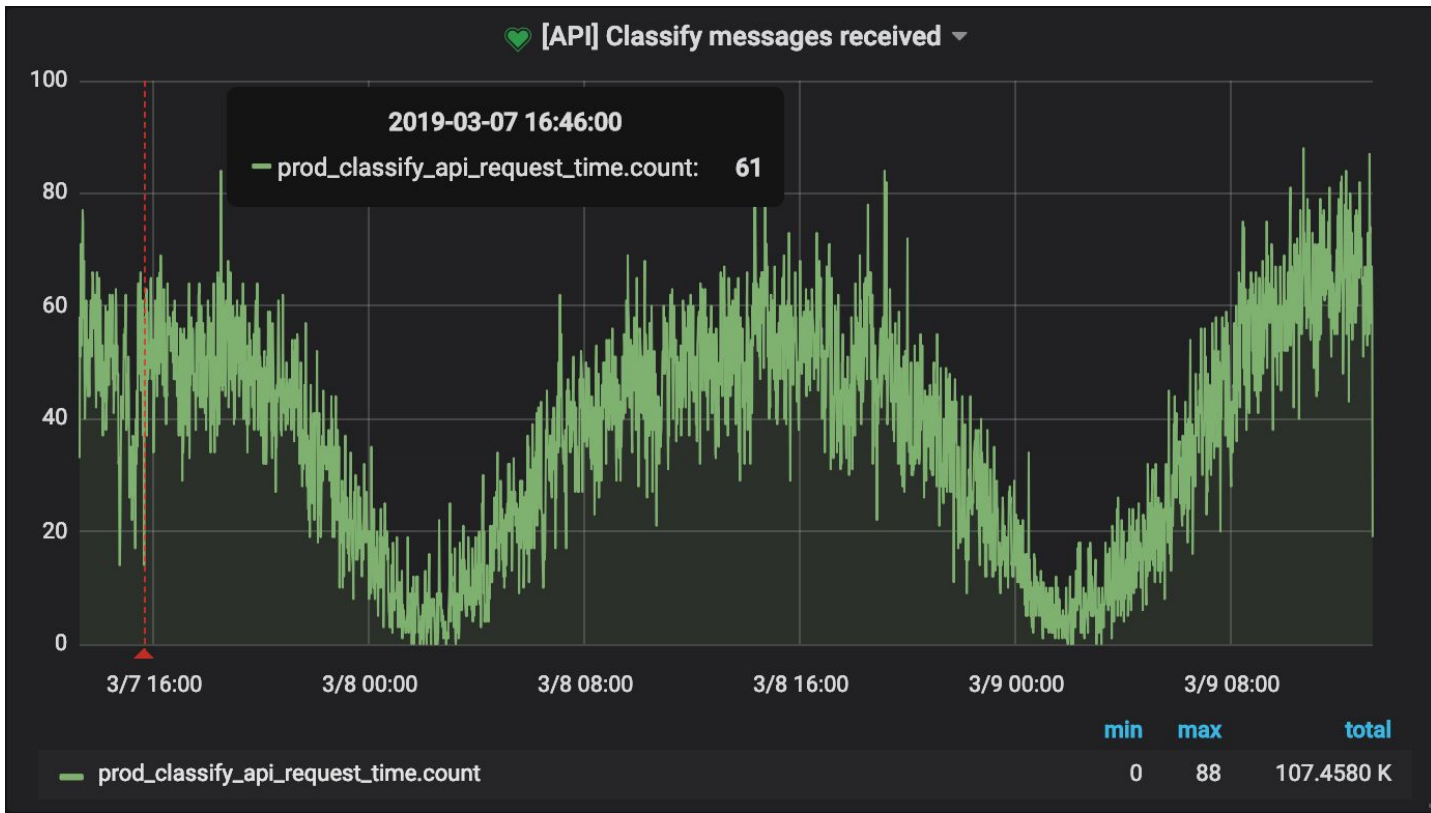


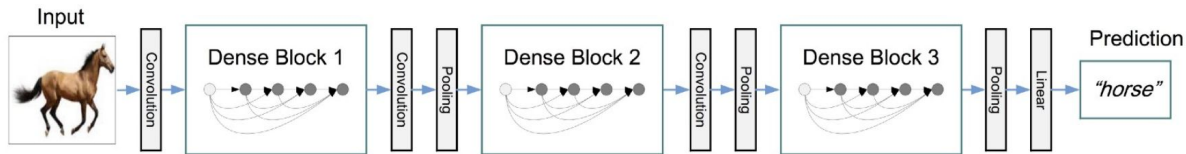
```
untitled
1 class ClassifyHandler(tornado.web.RequestHandler):
2     executor = ThreadPoolExecutor(2)
3
4     @run_on_executor
5     def predict(self, data):
6         results = self.classifier.predict(data)
7         self.storage: RedisStorage = self.get_storage()
8         self.storage.save(results)
9
10    async def post(self):
11        # get data
12        request_data = self.get_request_data()
13        await self.predict(request_data)
14
15    class GetClassifyHandler(tornado.web.RequestHandler):
16        async def get(self):
17            ad_id = self.get_request_data()['ad_id']
18            self.storage: RedisStorage = self.get_storage()
19            results = self.storage.get(ad_id)
20            self.finish({"status": "Success", "results": results})
21
22    application = tornado.web.Application([
23        (r"/classify", ClassifyHandler),
24        (r"/classify/results", GetClassifyHandler),
25    ])
```

StackOverflow: Ready to rescue



200 RPM

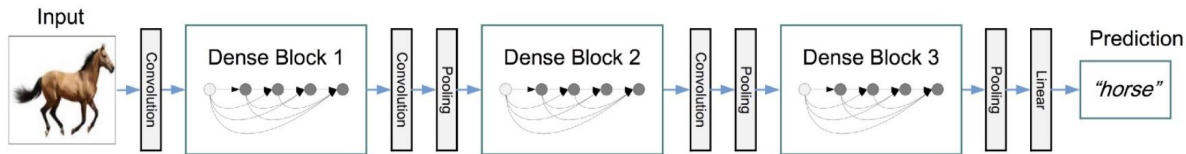




# PYTORCH







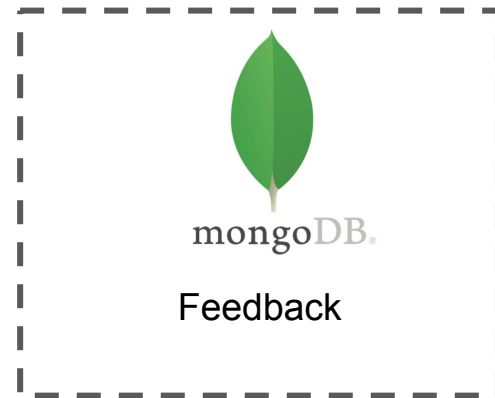
# PYTORCH



RAM - expiration?



Relational DB?



# Any issues there?

1. Request monitoring (?)
2. Hard to reason GPU usage (varies 3 to 7GB)
3. ~5Gb RAM out of the box (scale?)
4. Redis is in-memory storage (temporary)
5. Threading != Parallel



### PRESS IN DIRE SITUATIONS

If you're reading this on and Android, we've released the Nooooooooooooooooo button as a [Free App](#)



<http://www.noooooooooooooooooooo.com/>

# What if we split API and ML

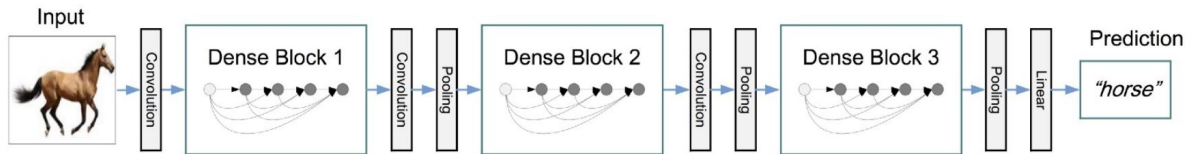


Part 2. A new hope

# Requirements

1. Single image/request processing time: < 2 seconds
2. Visibility
3. Persistence
4. Scalability
5. Make it extendable for new features
  - a. Price prediction
  - b. Similarity search
  - c. Segmentation
6. SDK friendly (well documented, tested etc)

Lalafo case



# PYTORCH

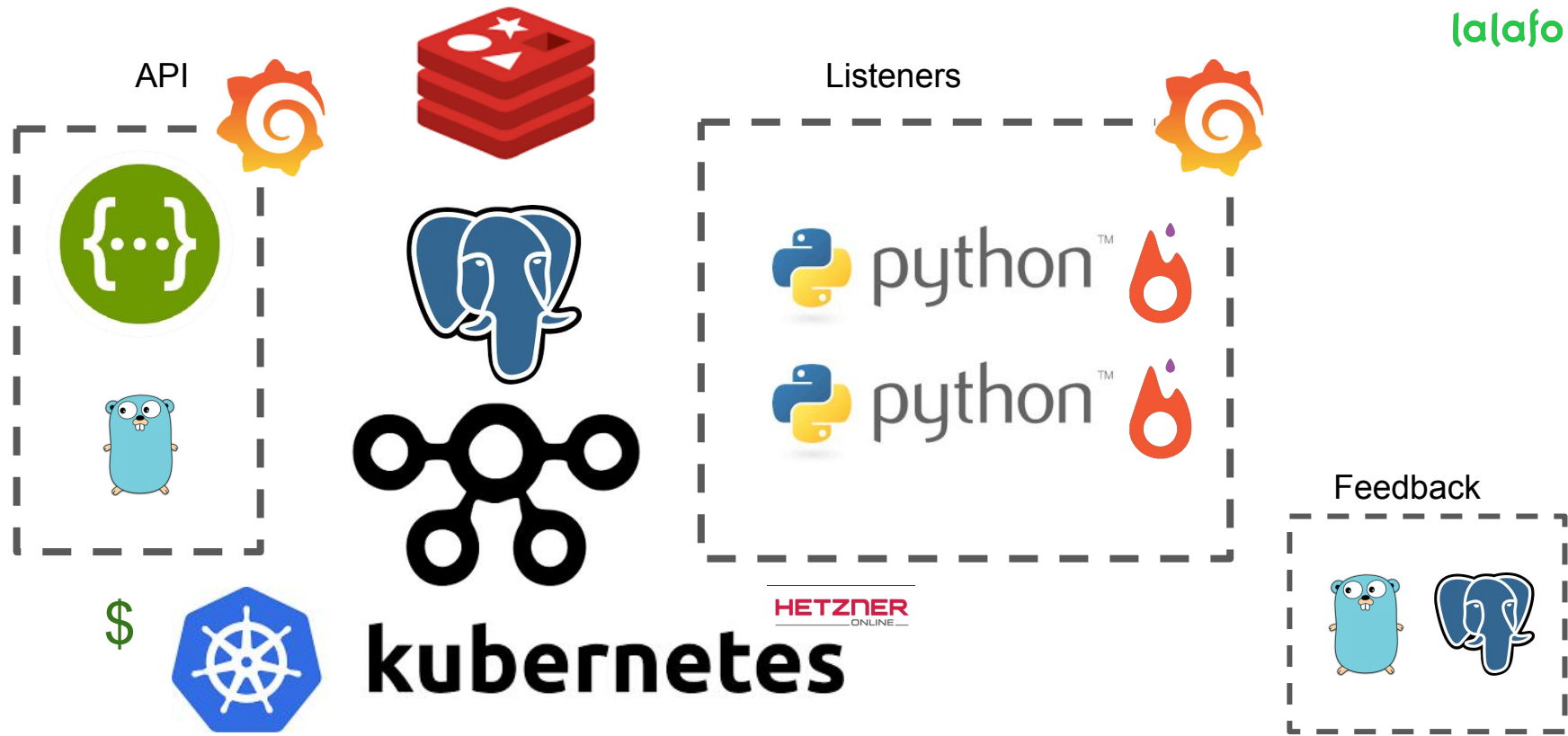


RAM - expiration?






Relational DB?





After 2 months

# Requirements

1. Single image/request processing time: < 2 seconds
2. **Visibility** - Decouple request and prediction 
3. **Persistence** 
4. **Scalability**  **kubernetes**
5. Make it **extendable for new features**
  - a. Price prediction
  - b. Similarity search
  - c. Segmentation
6. **SDK friendly** (well documented, tested etc)

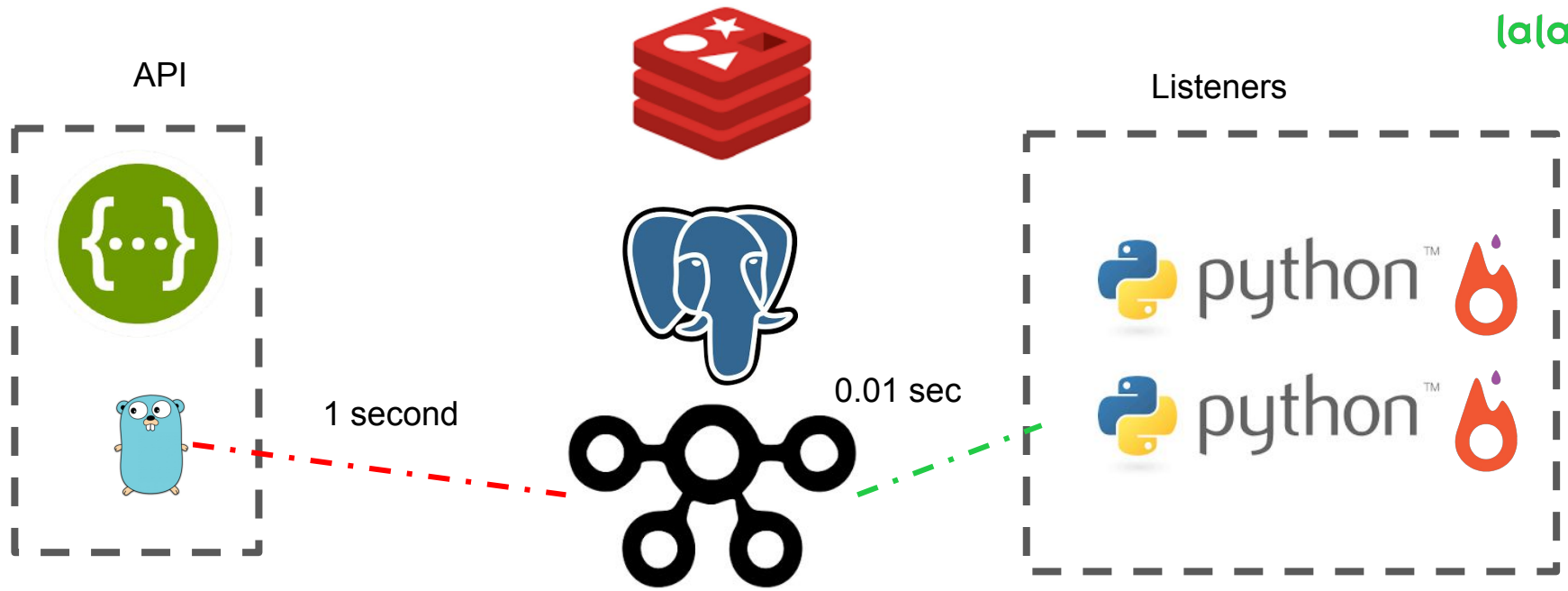
Lalafo case



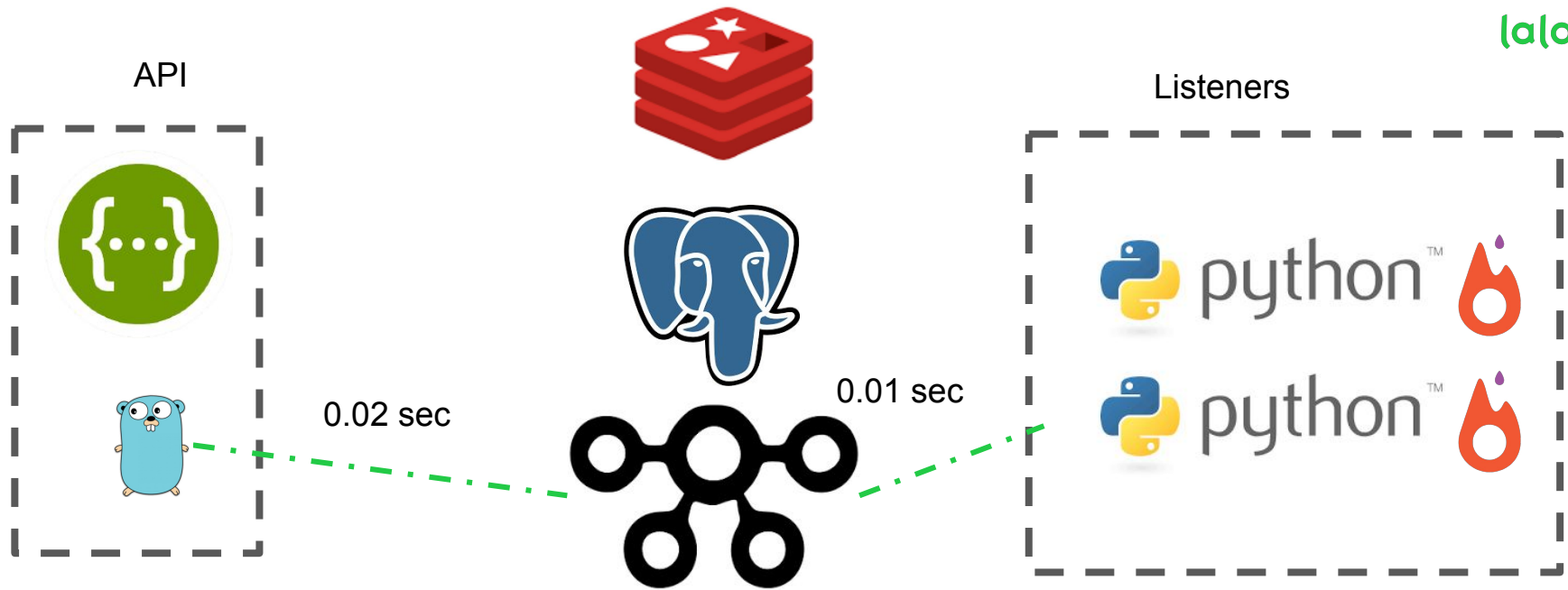
# 2 seconds per request



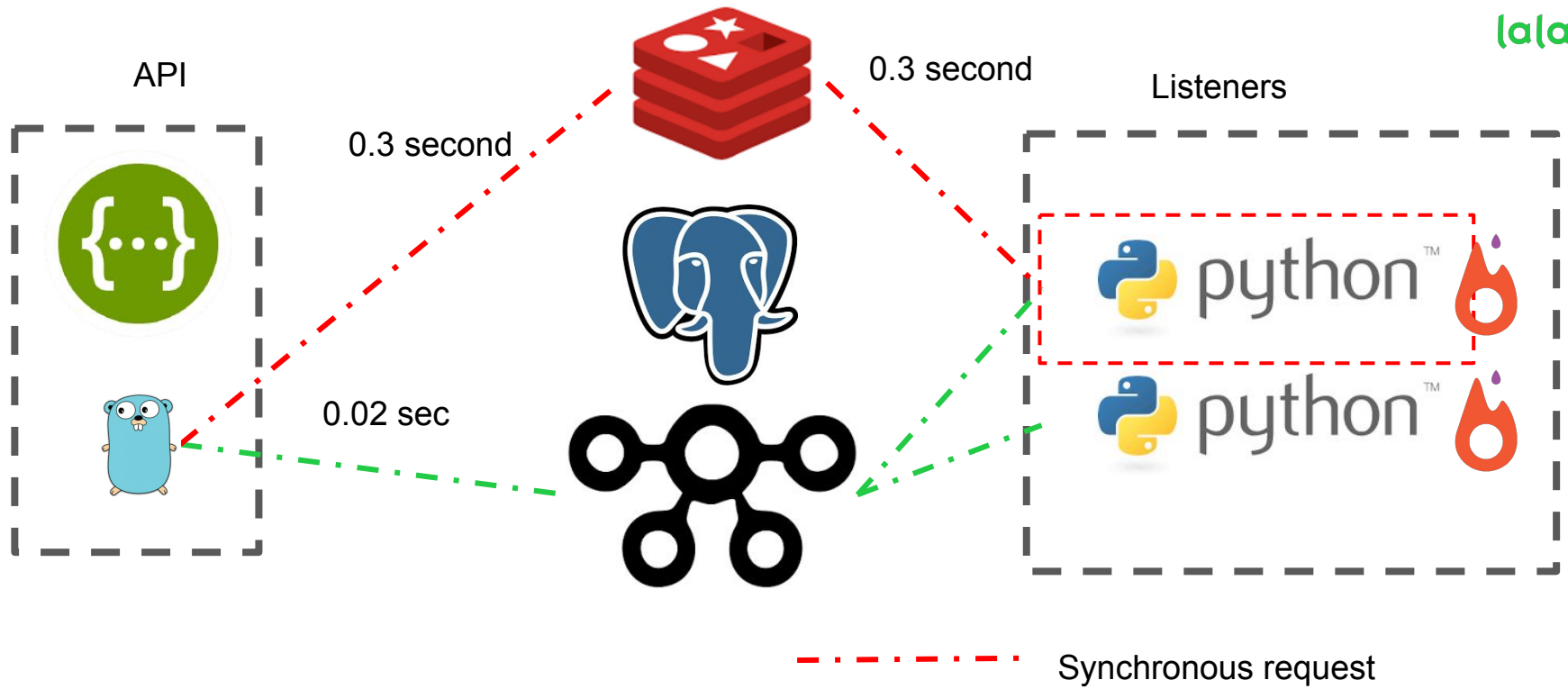
## Part 3. The Empire Strikes Back



4 seconds per request



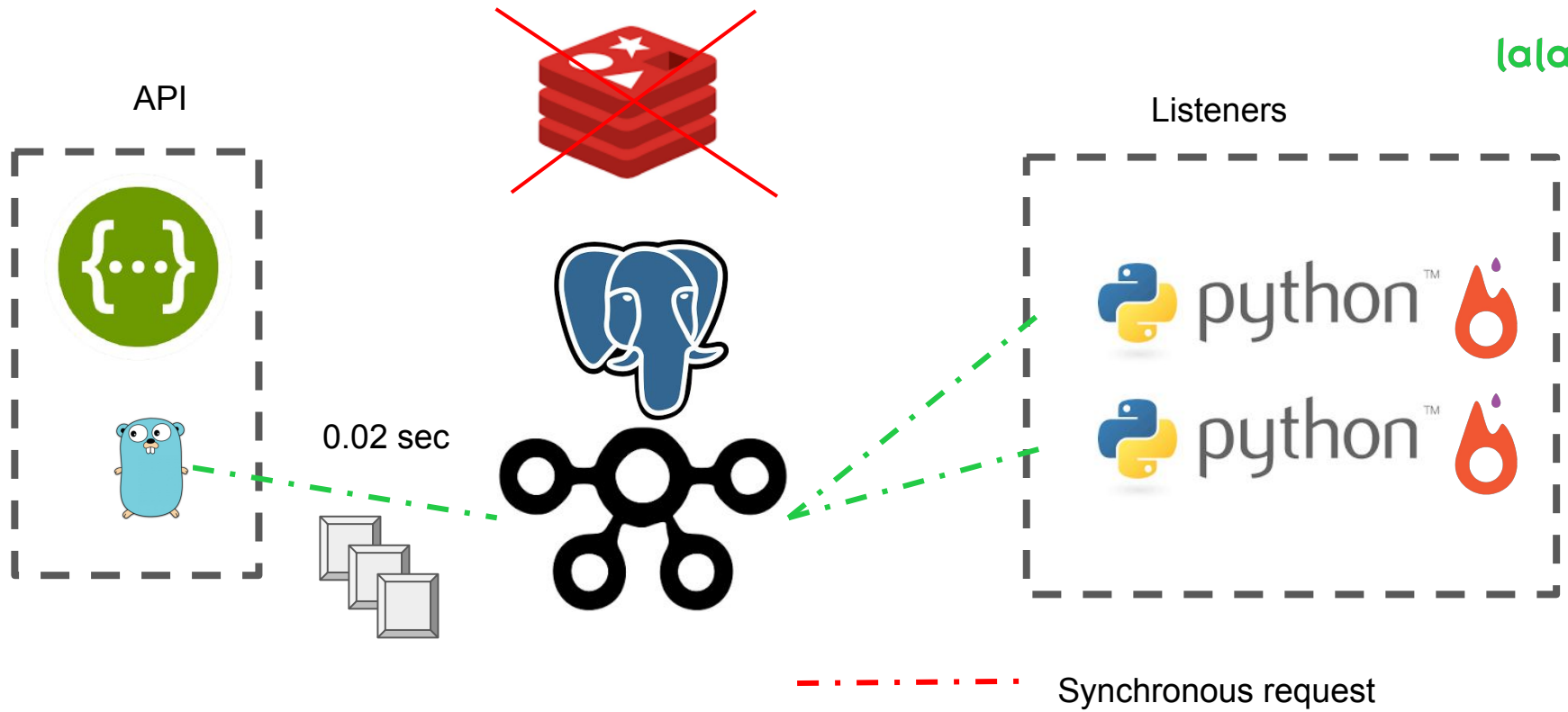
segmentio/kafka-go -> shopify/sarama



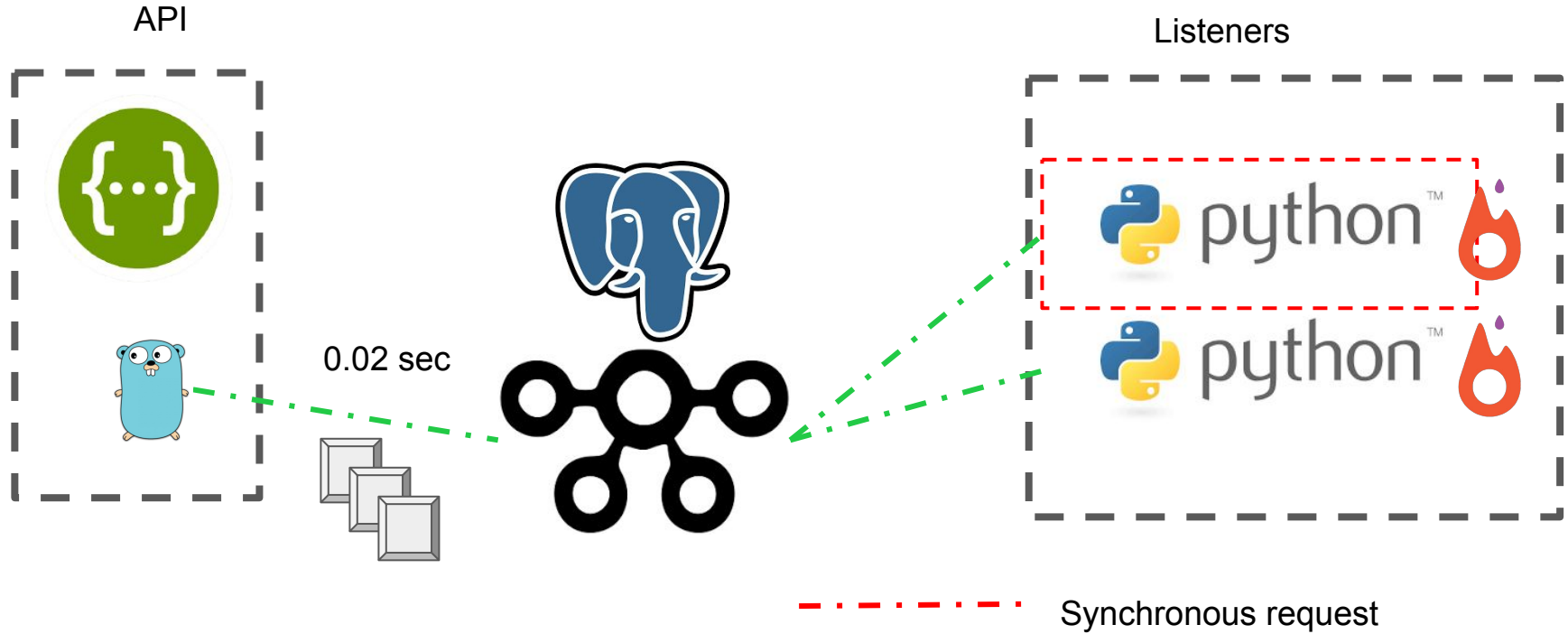
3 seconds per request



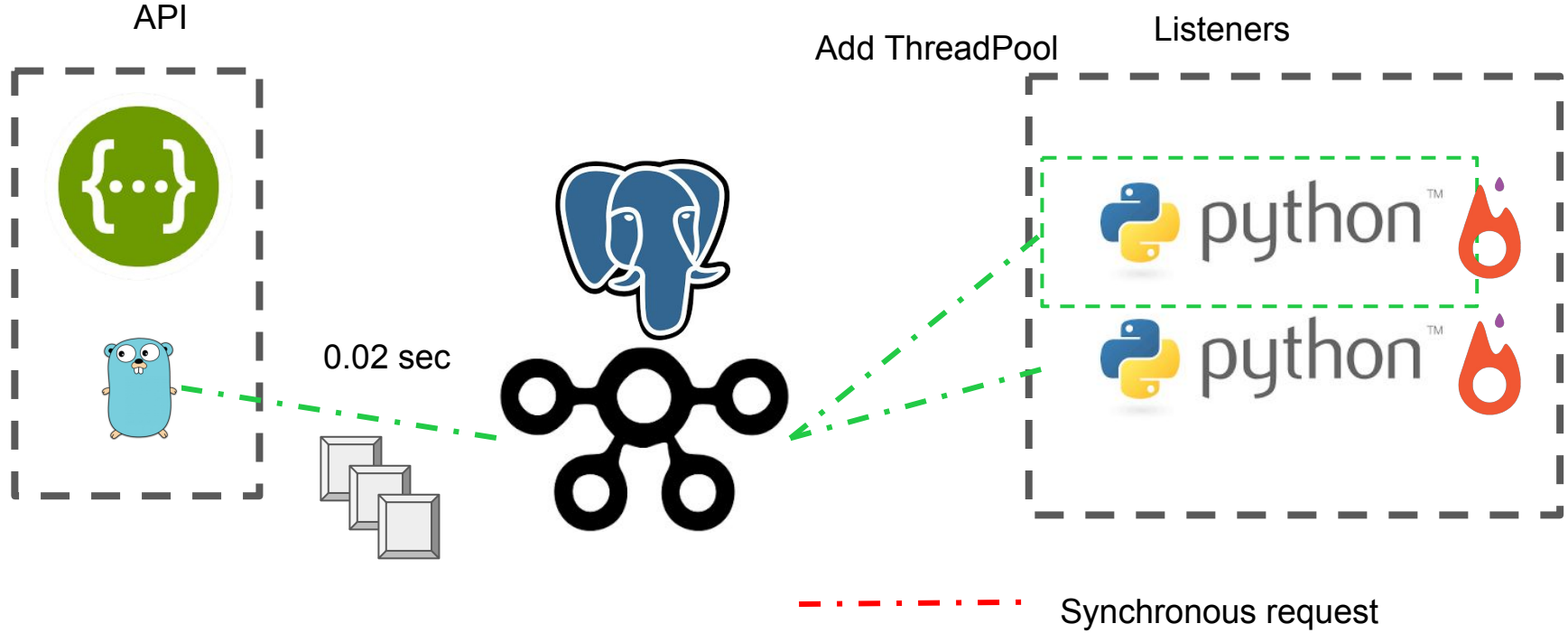
# How large is every image



2.4 seconds per request



2.4 seconds per request

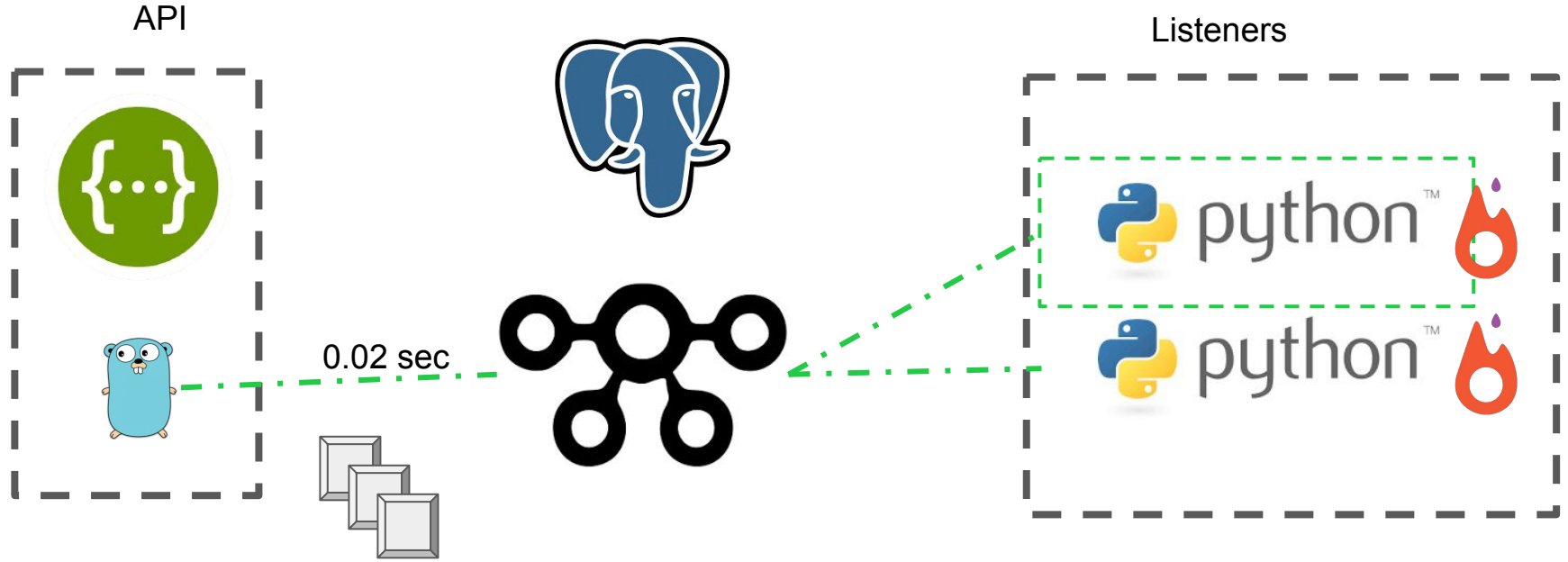


2 seconds per request

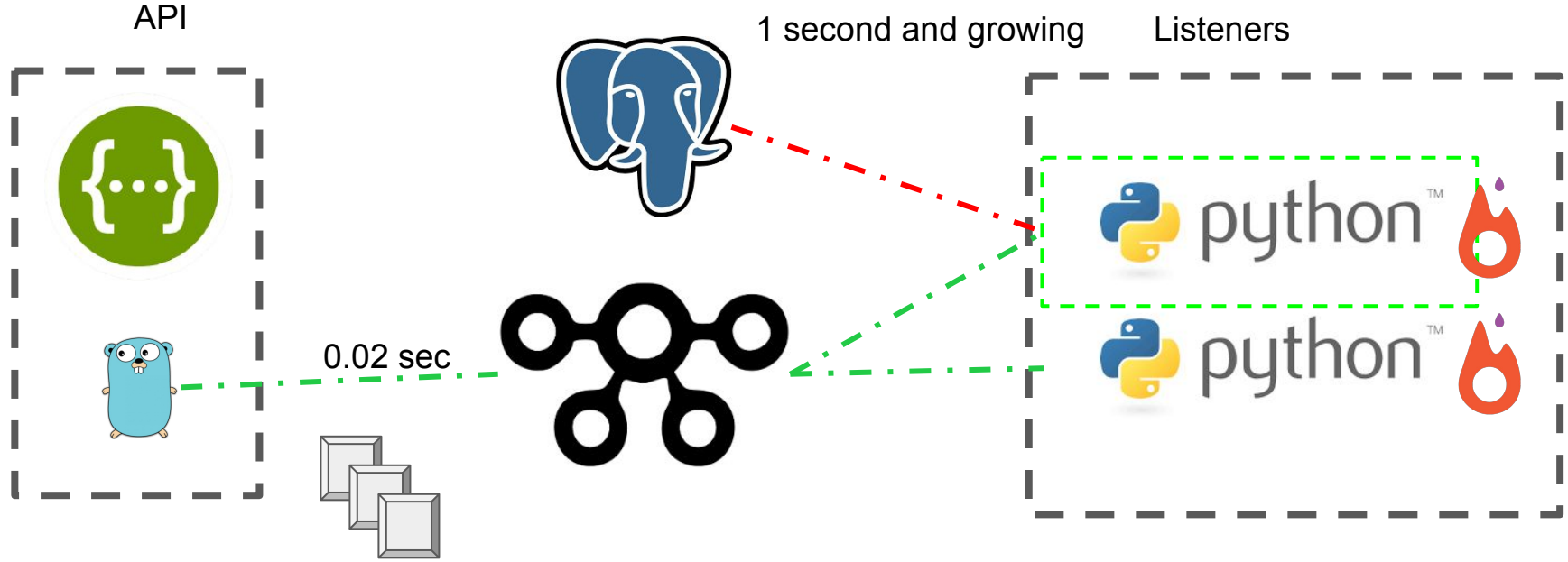


Success?  
Not yet

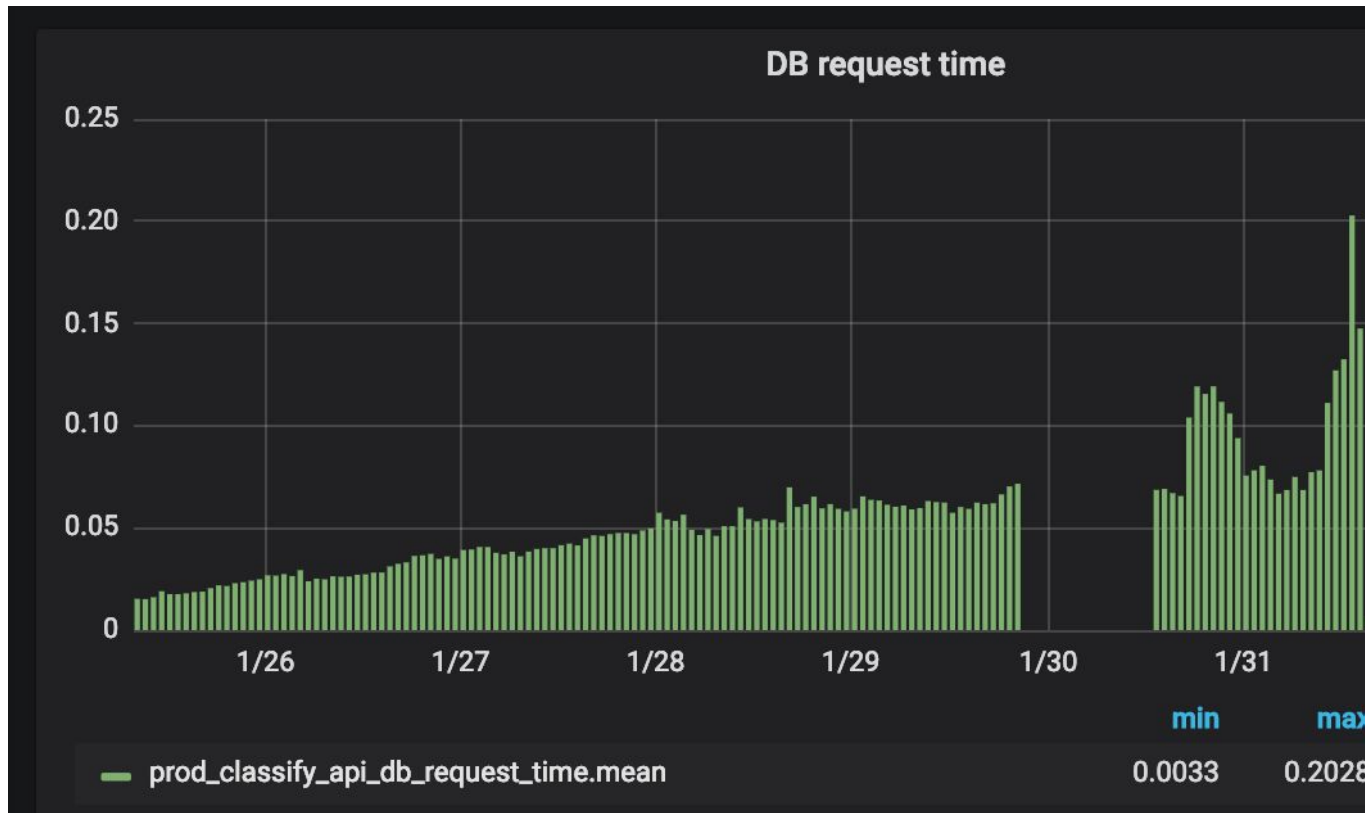
2 seconds per request



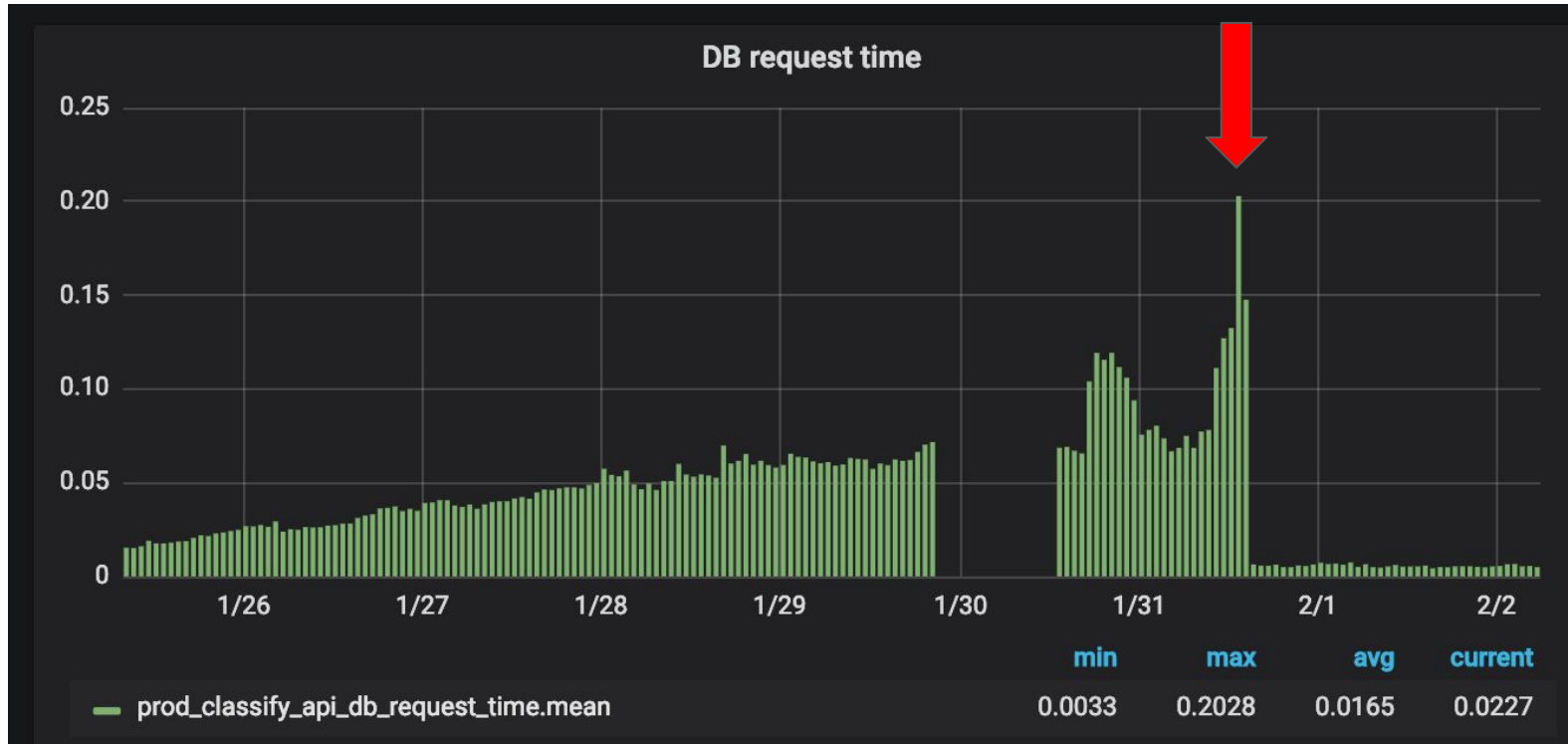
**3 seconds per request**



**3 seconds per request**



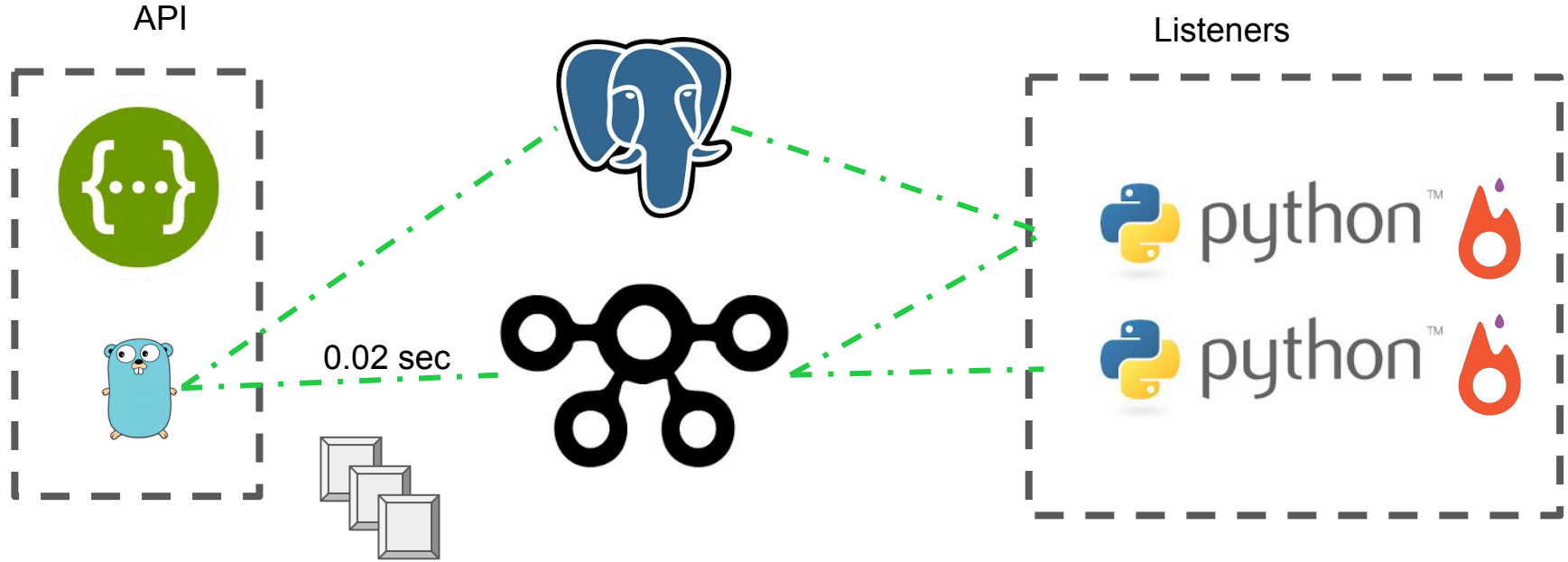
**3 seconds per request**



2 seconds per request



^-\_(ツ)^-

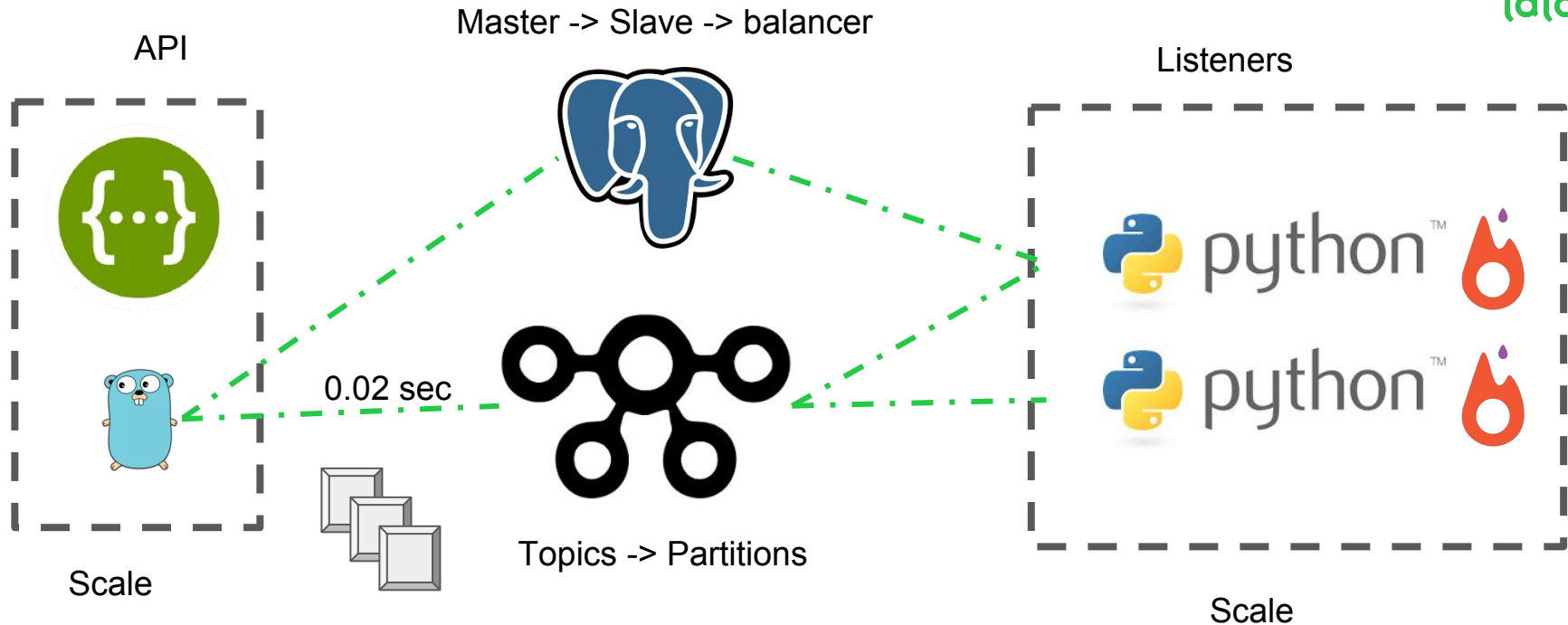


2 seconds per request



200 RPM -> 1200 RPM

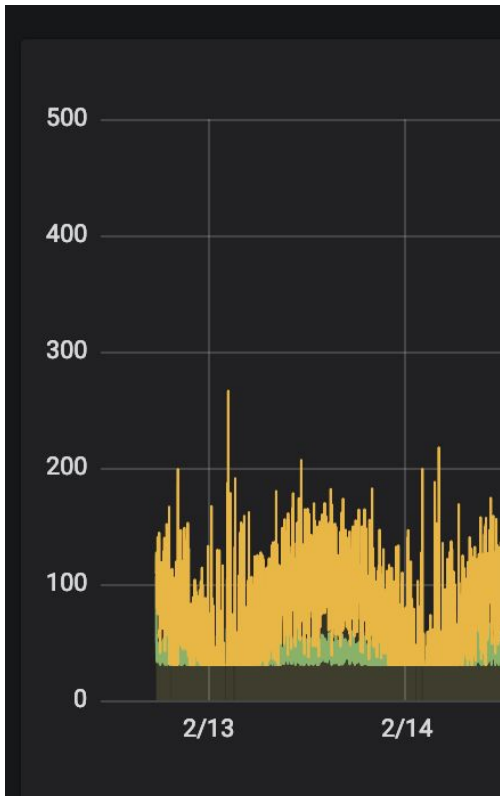




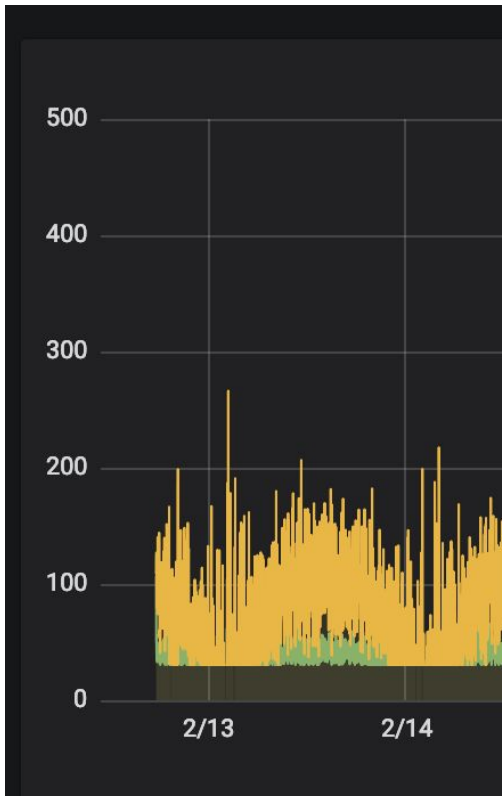
2 seconds per request

Success?  
Not yet

< 1 second per request



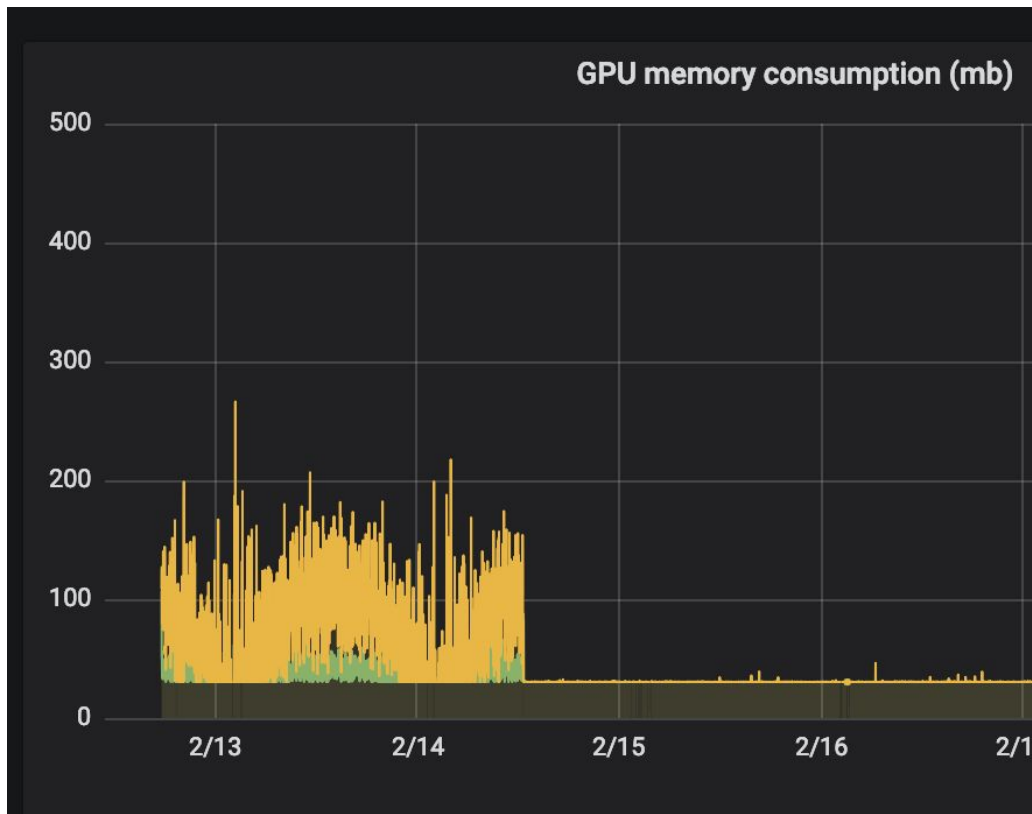
Processes:					GPU Memory
GPU	PID	Type	Process name		Usage
0	46897	C	/bin/python		7652MiB



Out of 8Gb



Processes:				GPU Memory Usage
GPU	PID	Type	Process name	
0	46897	C	/bin/python	7652MiB



How to PyTorch in production

Sat Mar 9 11:43:08 2019

```

+-----+
| NVIDIA-SMI 418.39      Driver Version: 418.39      CUDA Version: 10.1      |
+-----+-----+-----+
| GPU   Name           Persistence-M| Bus-Id        Disp.A | Volatile Uncorr. ECC |
| Fan  Temp  Perf    Pwr:Usage/Cap|      Memory-Usage | GPU-Util  Compute M. |
|=====+=====+=====+
|    0  GeForce GTX 1080    Off   | 00000000:01:00.0 Off  |          N/A   |
| 46%   57C    P2      40W / 180W | 1294MiB / 8119MiB |      0%      Default |
+-----+-----+-----+

```

```

+-----+
| Processes:                                     GPU Memory |
|  GPU       PID    Type    Process name                     Usage      |
|=====+=====+
|    0      15795     C   ...it- classify -p 6066 /app/classify-data  689MiB |
|    0     1546563     C      python3                      595MiB |
+-----+

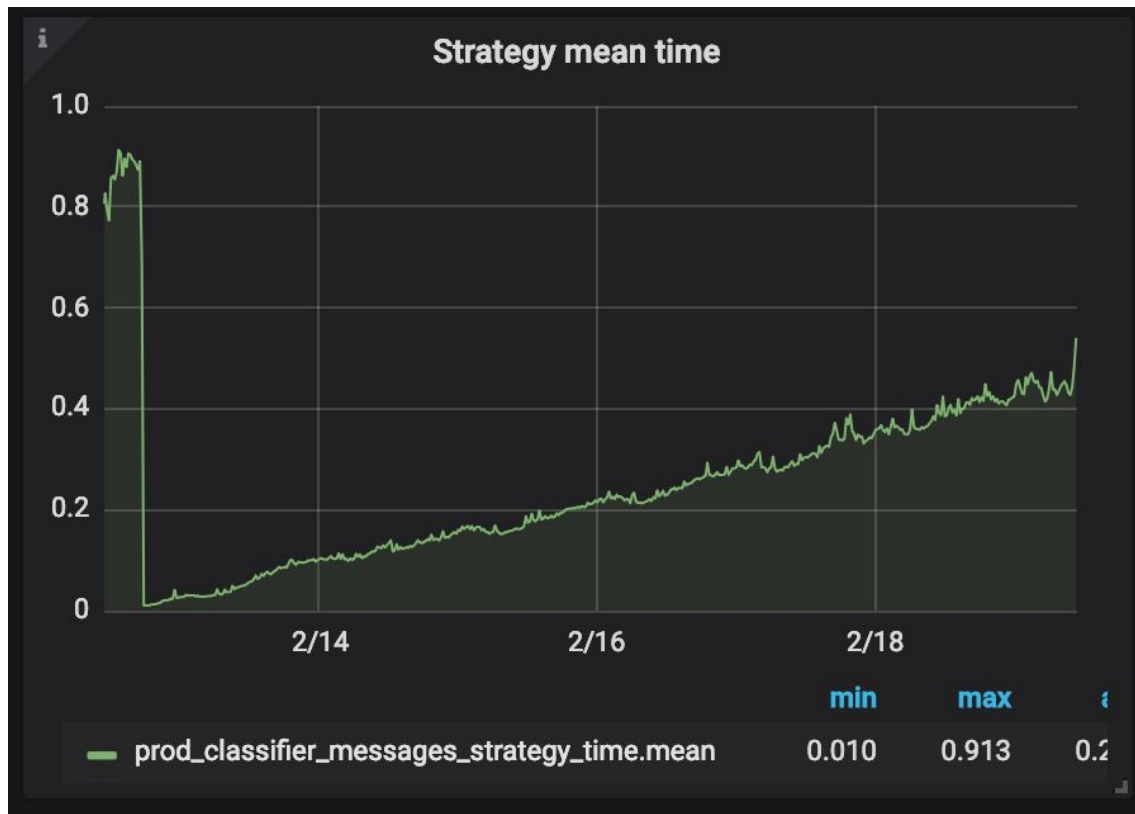
```



7.6 Gb -> 600 Mb

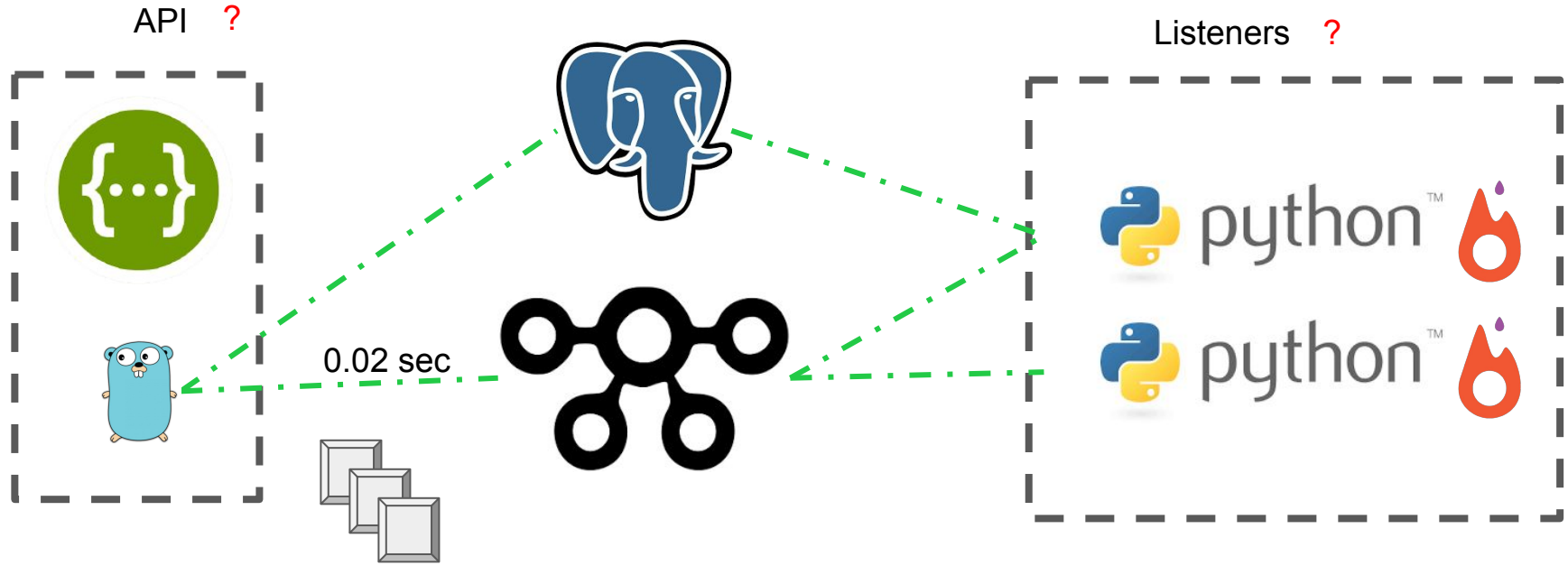
Success?  
Not yet

1 second per request

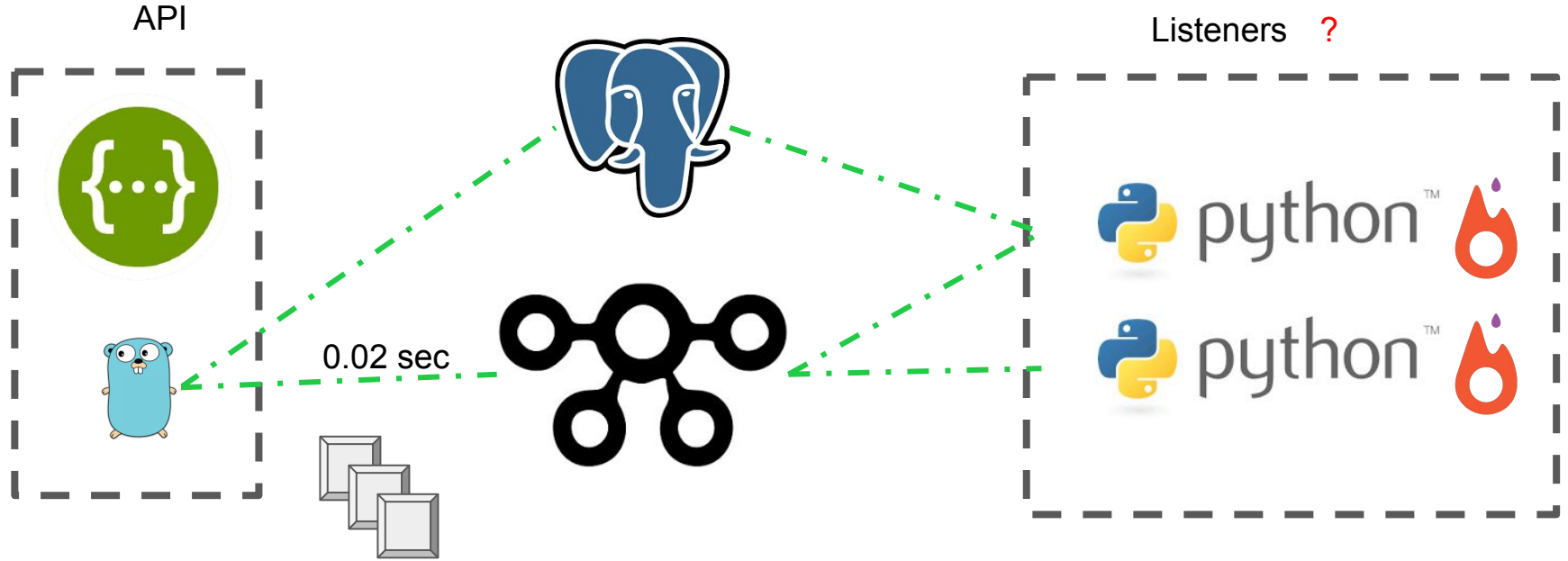


Back to 2, 5, 10 seconds per request

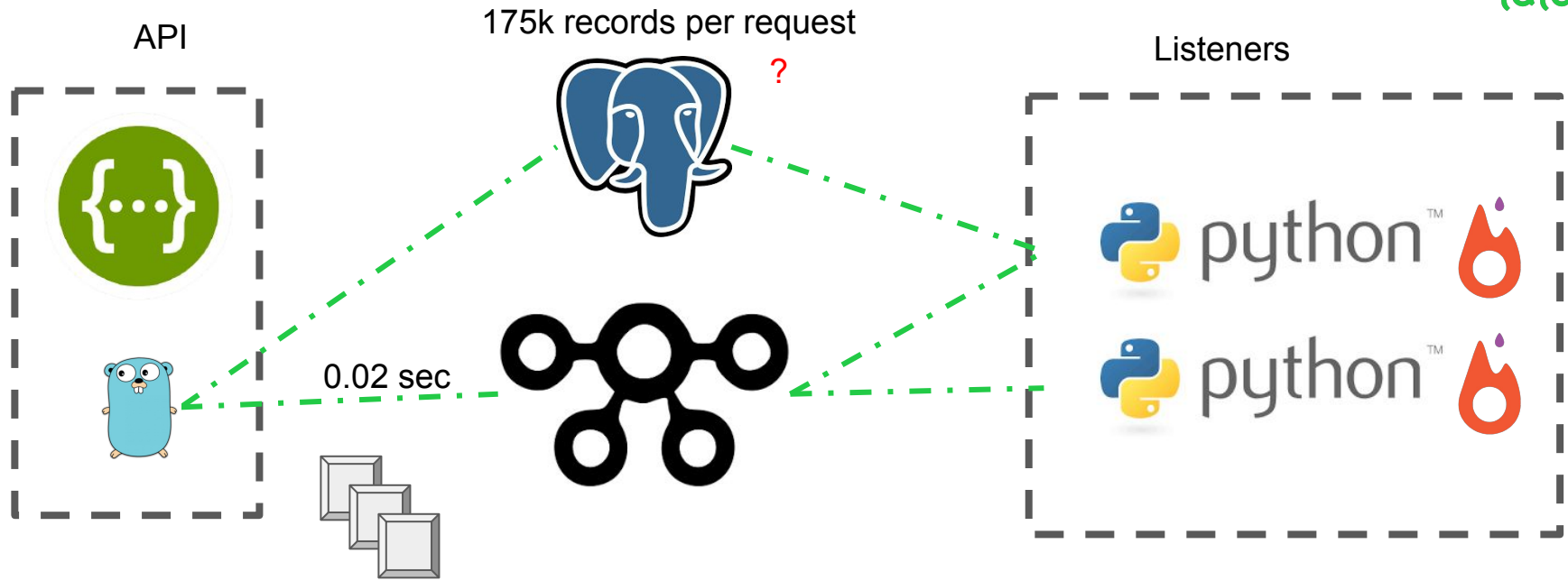




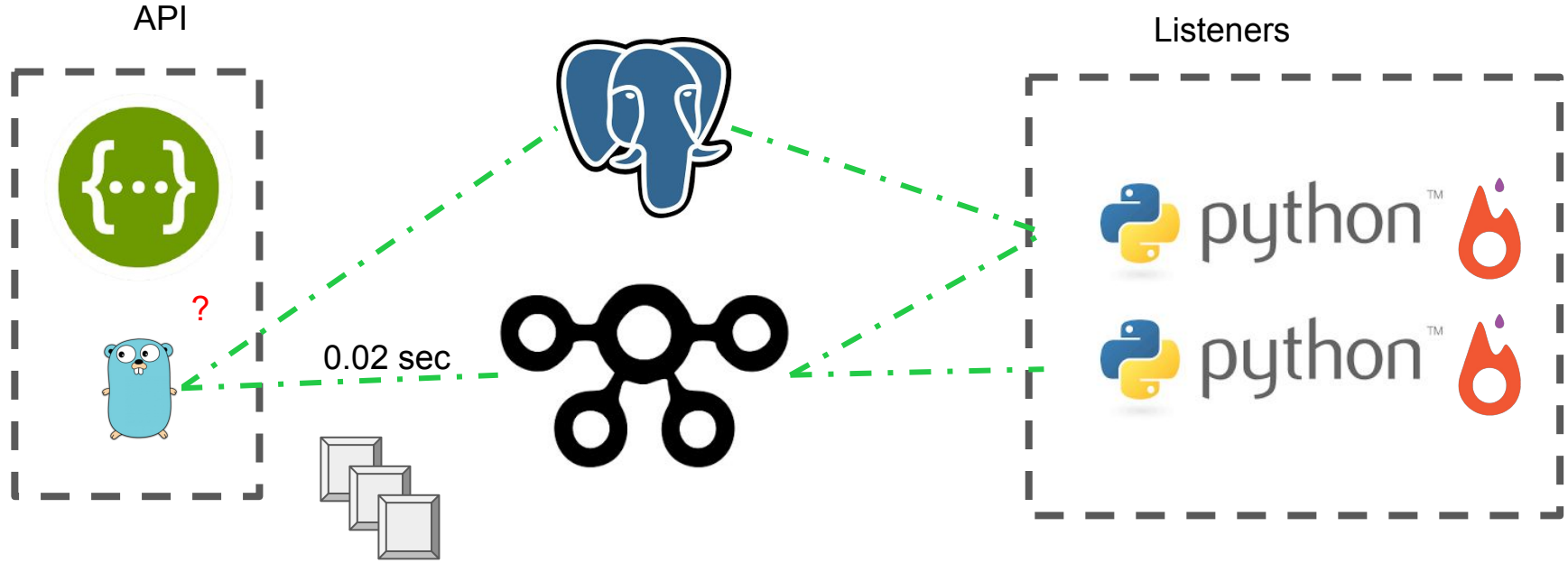
# Troubleshooting



# Troubleshooting



# Troubleshooting



# Troubleshooting



104	105		/
165		-	err := db.FirstOrCreate(&ad, Ad{CatalogAdID: ad.CatalogAdID, CatalogName: ad.CatalogName}).Error
	166	+	ad = Ad{CatalogAdID: ad.CatalogAdID, CatalogName: ad.CatalogName}
	167	+	err := db.Create(&ad).Error
166	168		influxTags := map[string]string{"action": "/v4/classify/add_images", "catalog": ad.CatalogName}
167	169		iw.WritePoint(dbRequestTime, influxTags, time.Since(funcStart).Seconds())
168	170		
169	171		if err != nil {
170		-	log.Errorf("couldn't save ad to the DB: %v\n Ad: %v\n", err, ad)
	172	+	if pqErr, _ := err.(*pq.Error); pqErr.Code == "23505" { // Retry once for integrity errors
	173	+	err = db.First(&ad).Error
	174	+	}
	175	+	log.Errorf("couldn't save ad to the DB: %v\n Ad: %v\n", err, ad)
171	176		c.JSON(http.StatusInternalServerError, LegacyApiError{Status: "Error", Message: "Could not

# Troubleshooting

104	105		
165		-	err := db.FirstOrCreate(&ad, Ad{CatalogAdID: ad.CatalogAdID, CatalogName: ad.CatalogName}).Error
	166	+	ad = Ad{CatalogAdID: ad.CatalogAdID, CatalogName: ad.CatalogName}
	167	+	err := db.Create(&ad).Error
166	168		influxTags := map[string]string{"action": "/v4/classify/add_images", "catalog": ad.CatalogName}
167	169		iw.WritePoint(dbRequestTime, influxTags, time.Since(funcStart).Seconds())
168	170		
169	171		if err != nil {
170		-	log.Errorf("couldn't save ad to the DB: %v\n Ad: %v\n", err, ad)
	172	+	if pqErr, _ := err.(*pq.Error); pqErr.Code == "23505" { // Retry once for integrity errors
	173	+	err = db.First(&ad).Error
	174	+	}
	175	+	log.Errorf("couldn't save ad to the DB: %v\n Ad: %v\n", err, ad)
171	176		c.JSON(http.StatusInternalServerError, LegacyApiError{Status: "Error", Message: "Could not

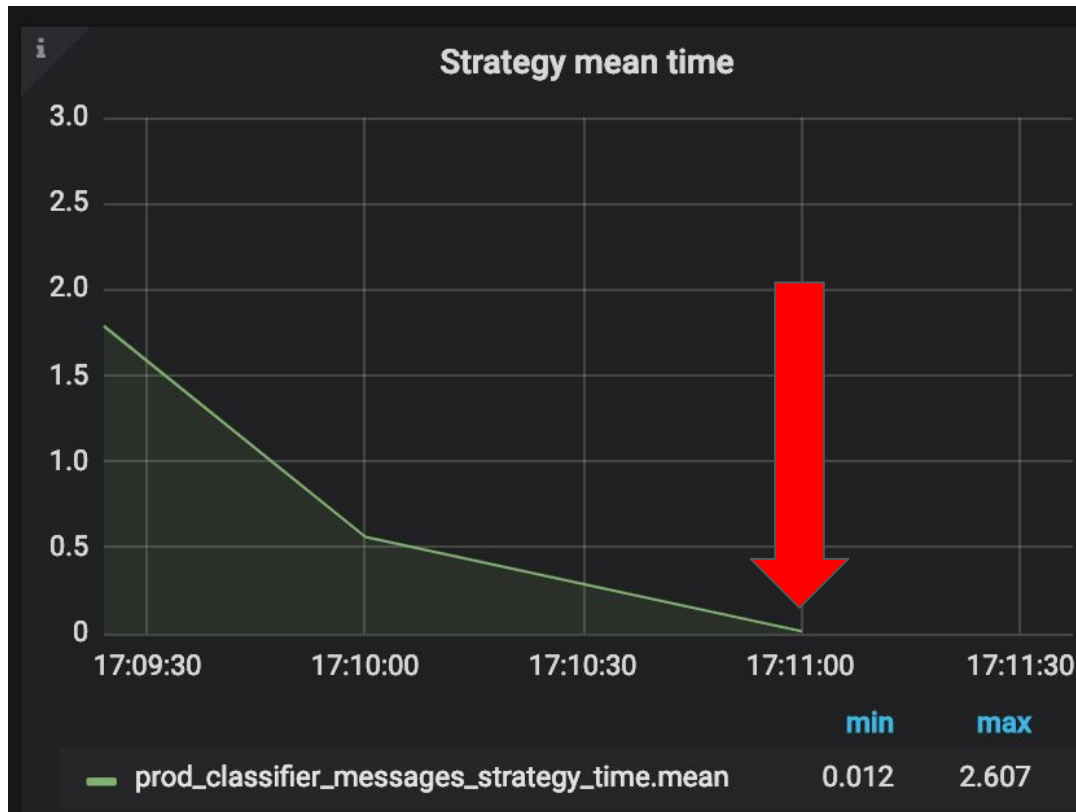


			"uix_catalog_name_ad_id\"" {
178		-	err = db.First(&ad).Error
	178	+	err = db.First(&ad, &Ad{CatalogName: ad.CatalogName, CatalogAdID: ad.CatalogAdID}).Error

# SELECT \* FROM ads LIMIT 1



25 seconds -> 0.02 seconds



0.02 seconds -> 0.002 second



Success?  
Not yet



**Ask Solem**  
@asksol

helping make distributed systems easy to build for everybody. Created Faust with vineetik and before that Celery. Works at @robinhoodapp 🌳🌳🌴🌵

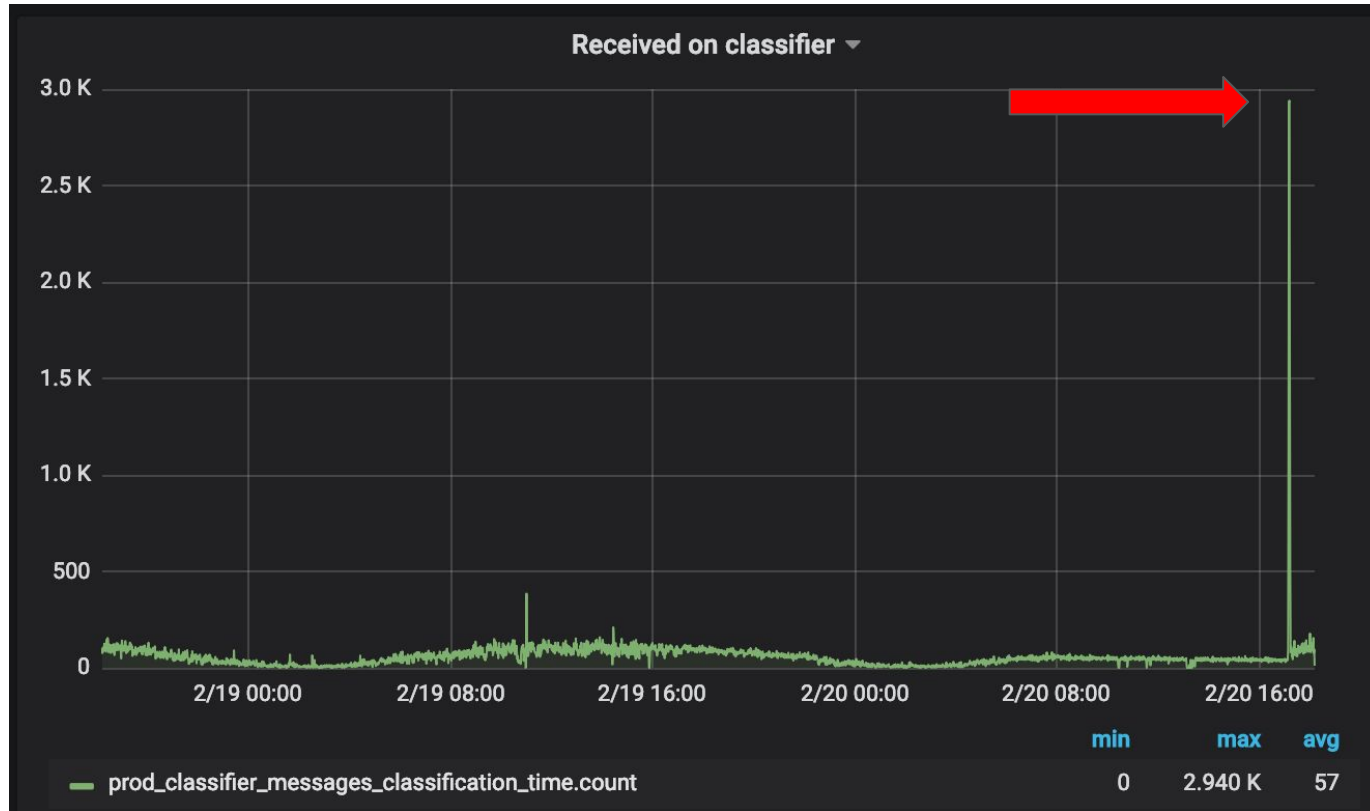
📍 San Francisco, CA

<https://github.com/robinhood/faust>



F A U S T

# Faust 1.4.6: No latest offset



3K images per minute 0\_o

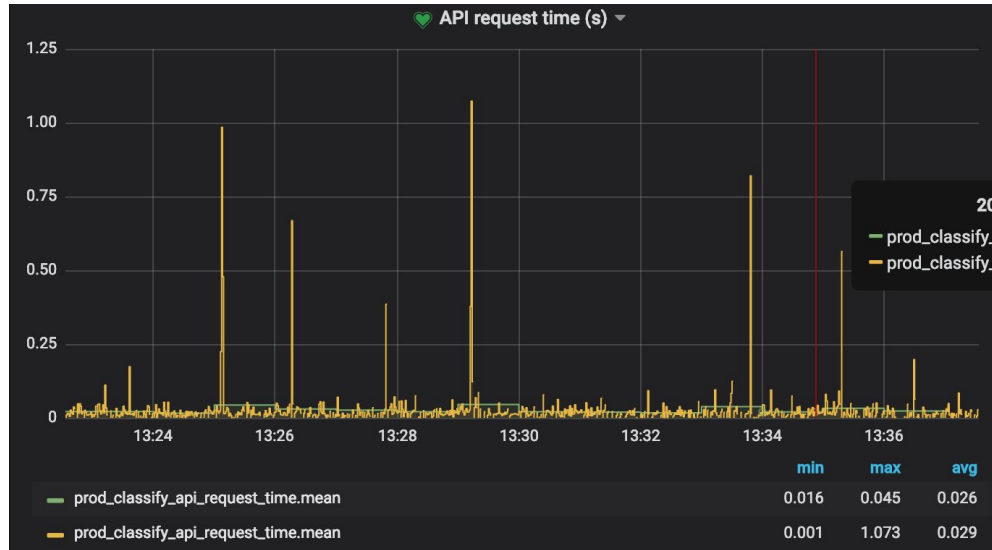
Success?  
Not yet

# Issues to solve

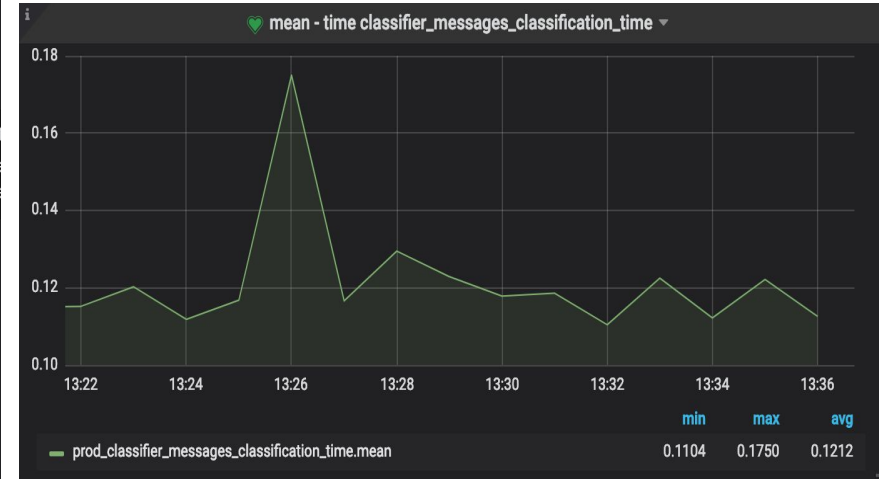
1. Occasional spikes in performance (GC, network latency)
2. Message broker (Kafka rebalancing, offset etc)
3. How to handle DB migrations
4. Something we are not aware of yet

# Lessons learnt

1. CPU bound tasks != IO bound (¯\\_(\ツ)\\_/\)
2. High coupling - low cohesion
3. You need to know how to cook MongoDB
4. Go is not that obvious and library reach as Python
5. Simple != Easier
6. Concurrency != Parallelism (obviously)

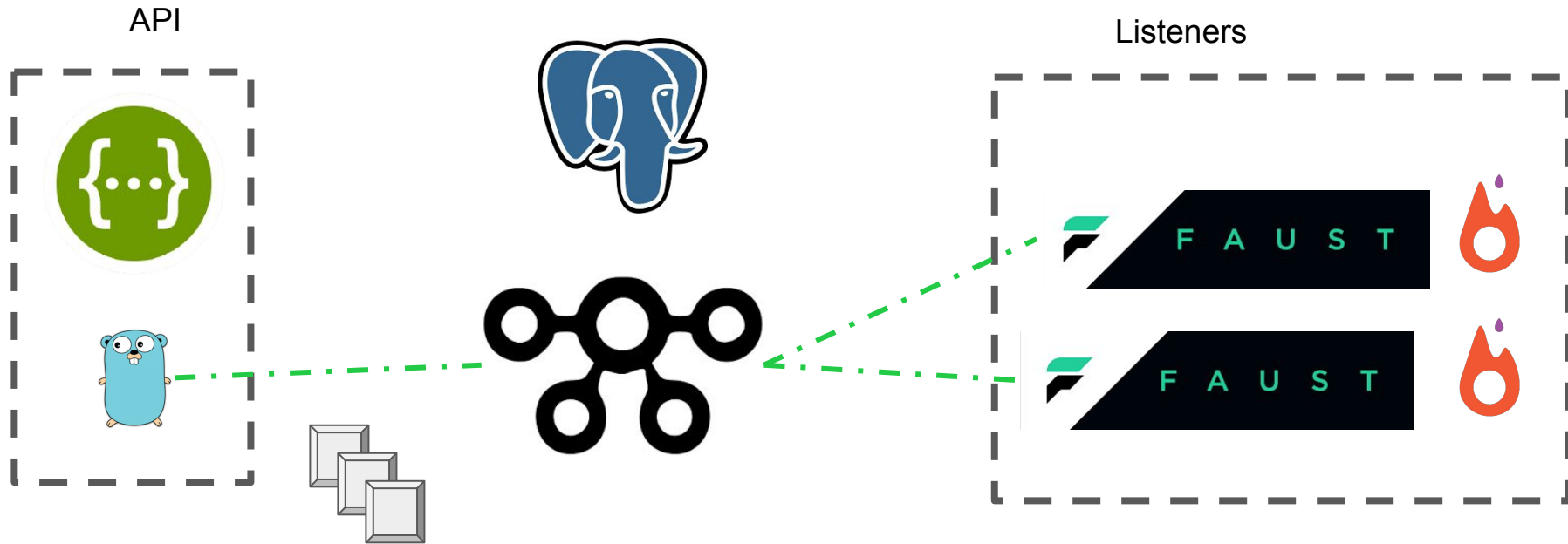


min	max	avg
0.015	0.045	0.027
0.001	1.073	0.029



min	max	avg
0.1025	0.1750	0.1208

< 1 second per request



+ Live statistics from PostgreSQL



# Success?

Thank you everyone