

The logo for Coveo. It features the word "coveo" in a lowercase, white, sans-serif font. A small, solid orange triangle is positioned above the letter "v". A trademark symbol (TM) is located at the top right of the letter "o".

coveo™



# Passer de 5,000 requêtes à 50,000,000 par jour

Récit des apprentissages



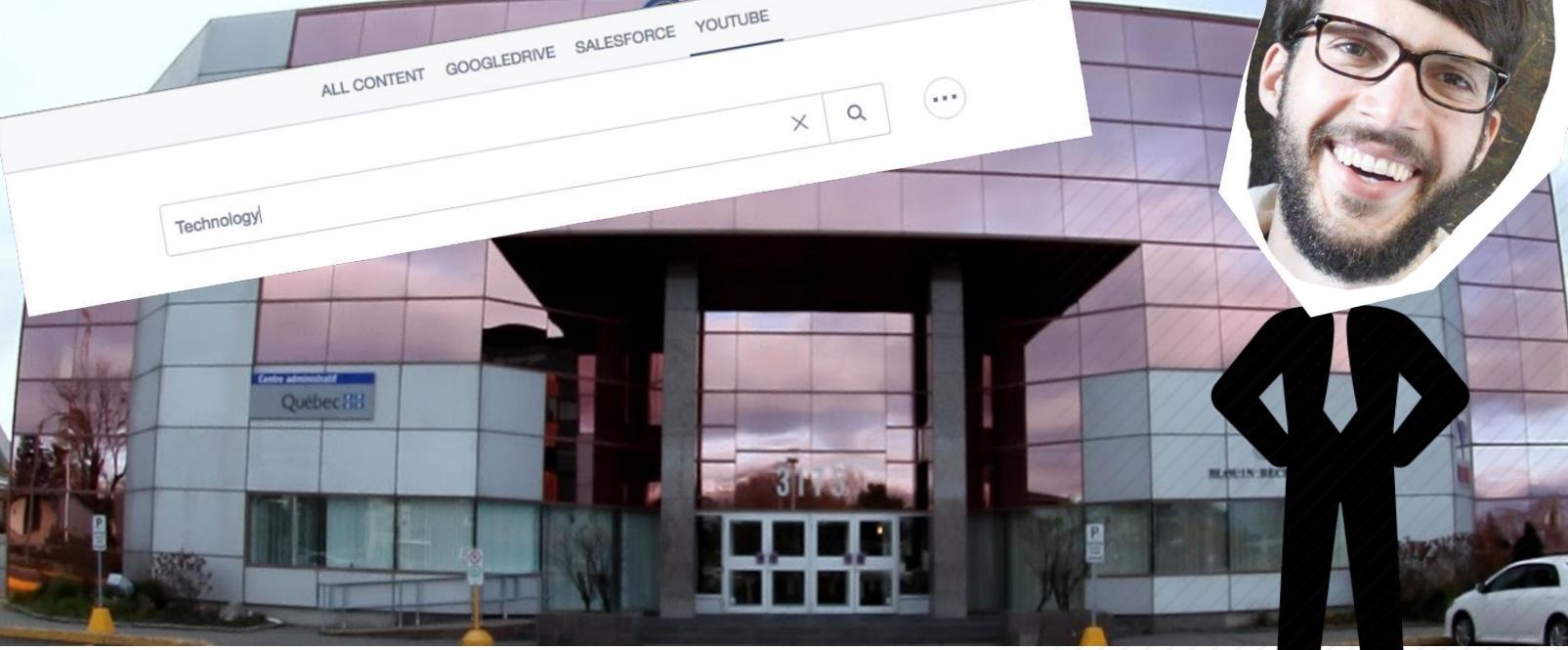
Andy Emond

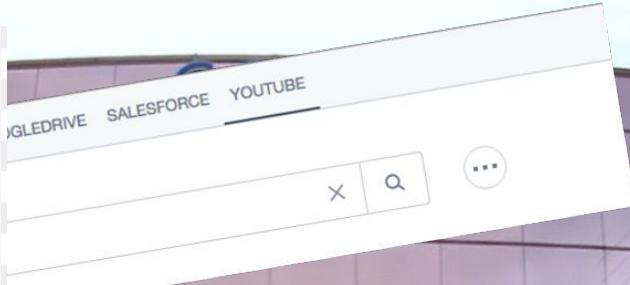
DEVELOPPEUR LOGICIEL, COVEO

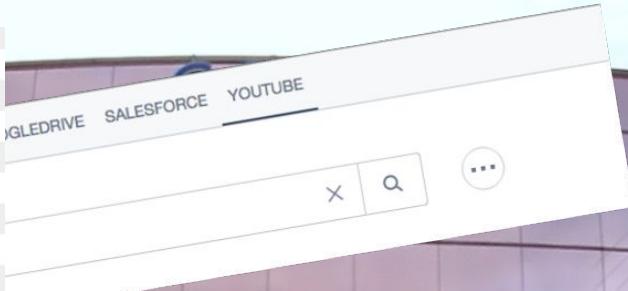
# Mise en contexte Coveo et MLBackend









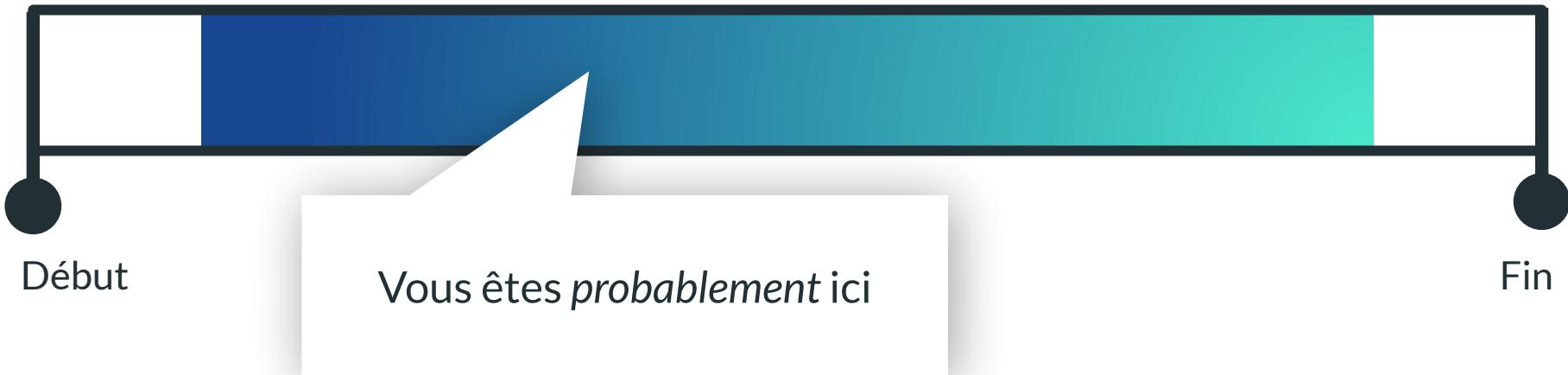




# Ligne de vie d'un projet



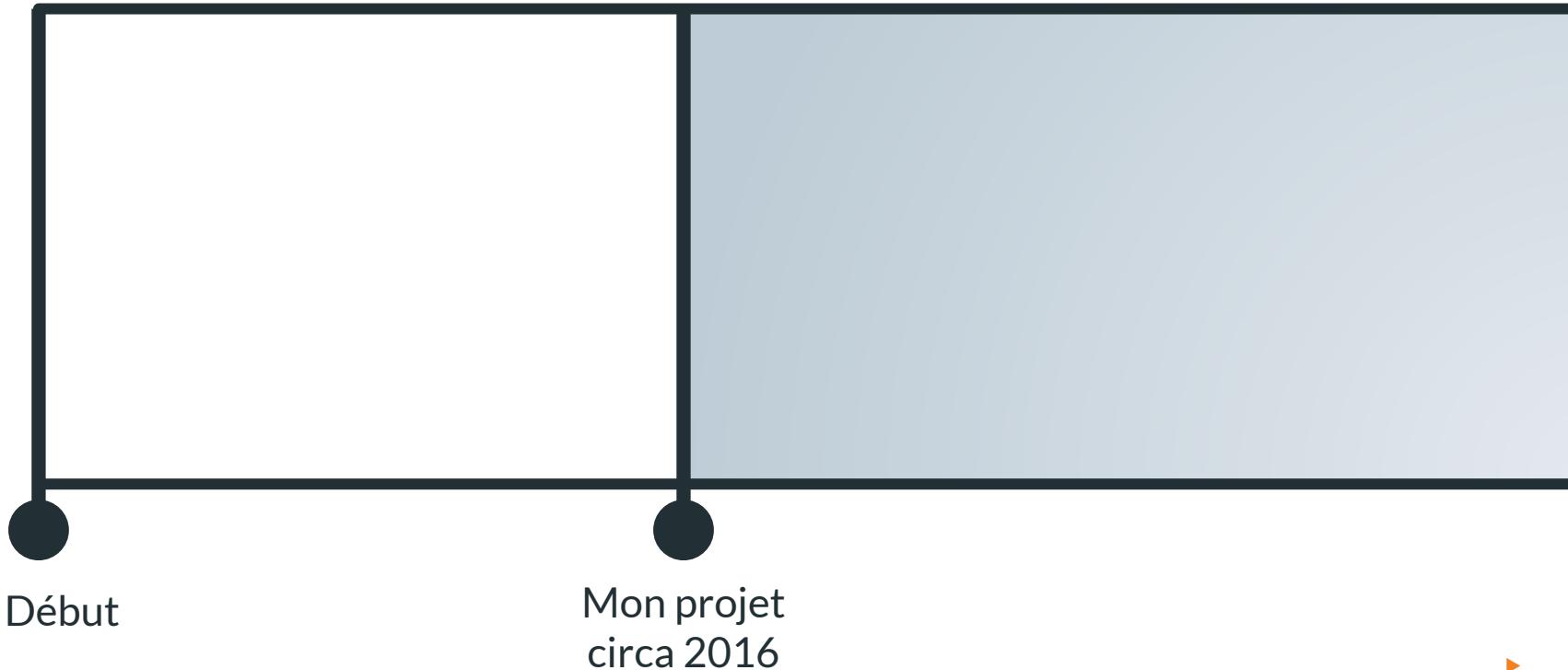
# Ligne de vie d'un projet



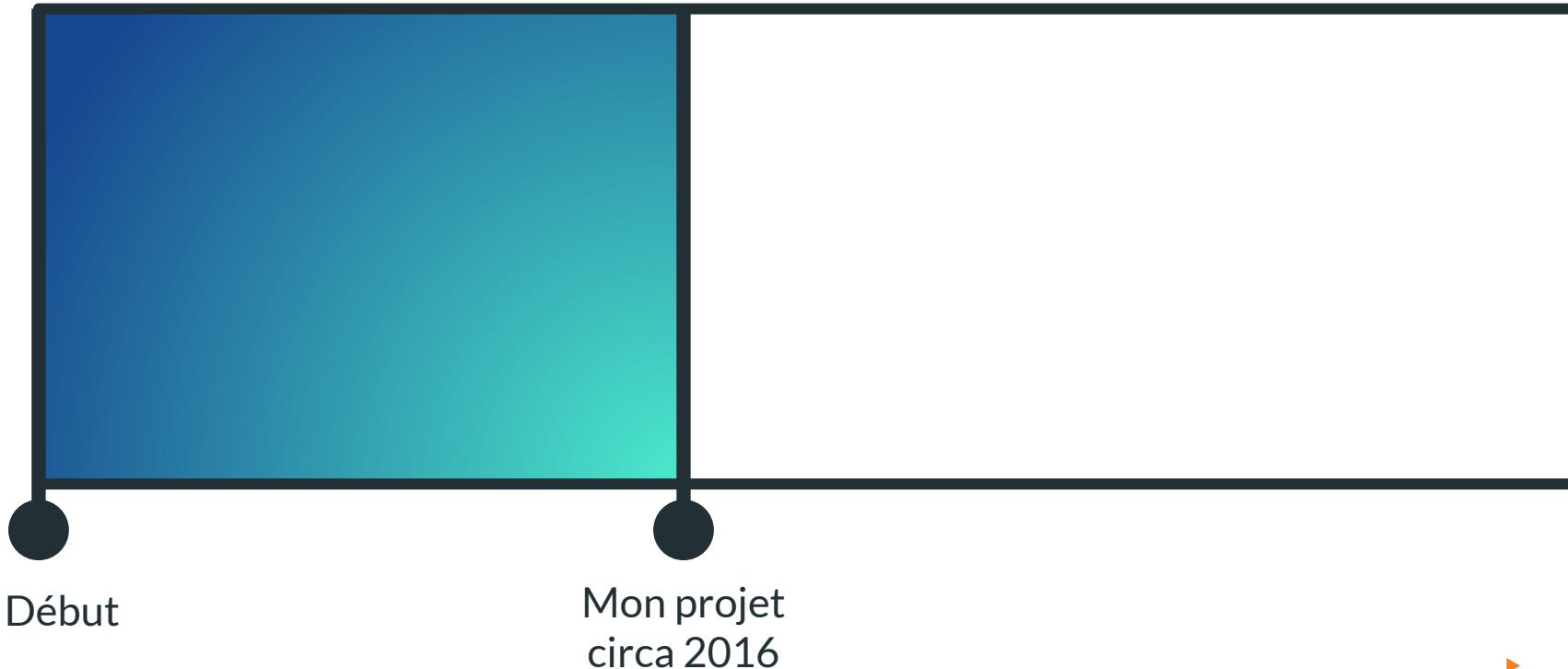
# Ligne de vie de **mon projet**



# Ligne de vie de **mon projet**



# Ligne de vie de **mon projet**



# Genèse du projet



Né d'un  
**hackathon**

Et a survécu!



# Make it **work**™

Taillé sur pièce pour le  
problème désiré



Est-il possible  
de faire du  
“ML”?

Oui.

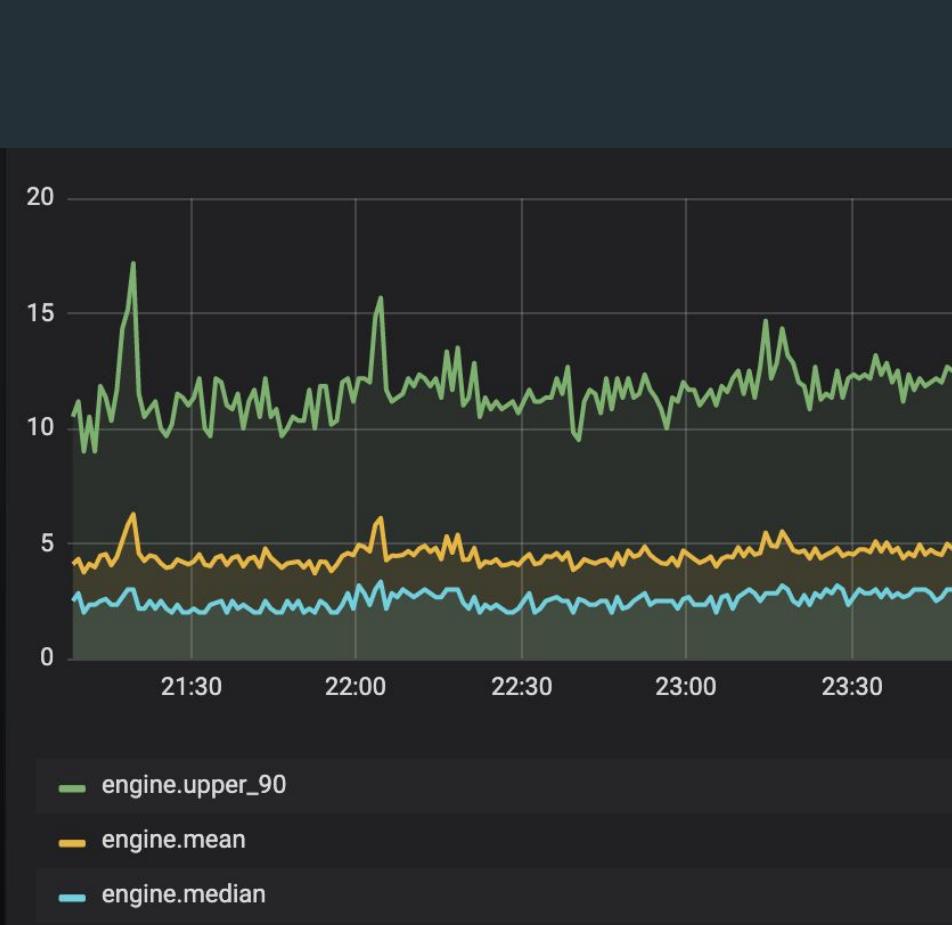
Est-il possible de faire  
du “ML” à grande échelle?

# La situation

# La situation?

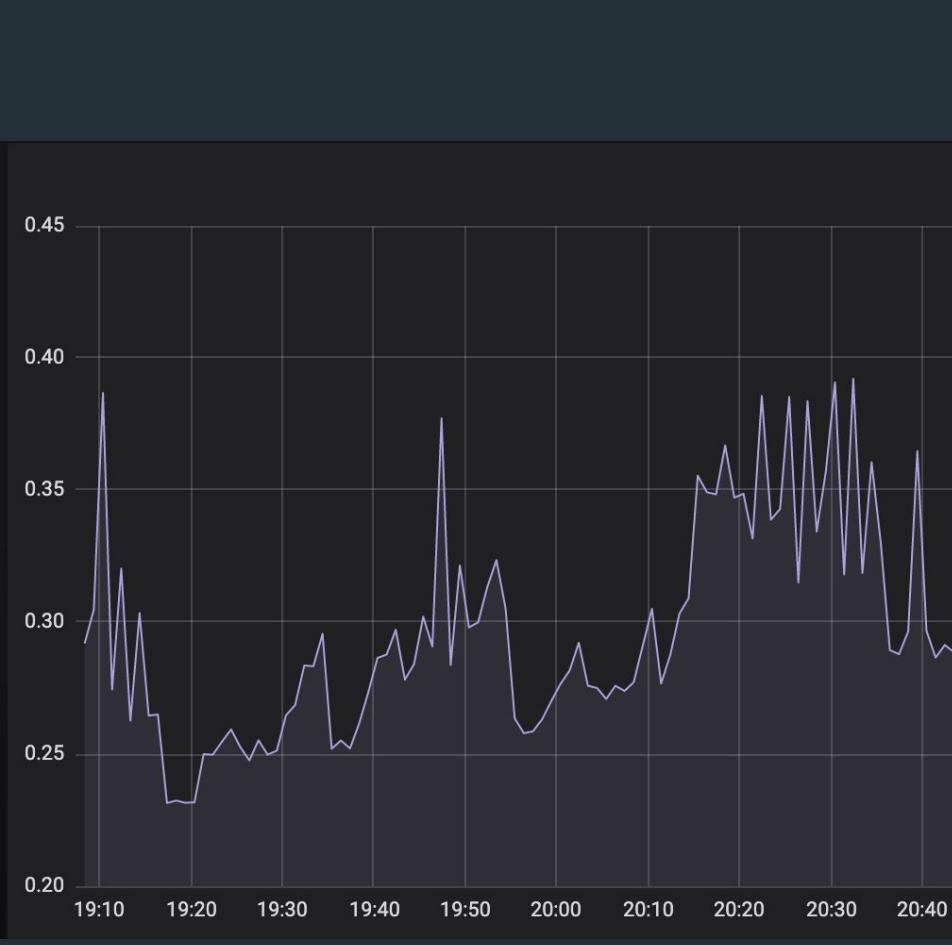
“S'il est impossible de le mesurer, il est impossible de l'améliorer.”

PETER DRUCKER



# Temps de réponse

## Percentile



## Temps de réponse Percentile

## Utilisation mémoire Pression sur le Garbage Collector



## Temps de réponse

### Percentile

## Utilisation mémoire

### Pression sur le Garbage Collector

## Utilisation CPU

Si une machine tombe, est-ce que les autres ont l'espace nécessaire pour soutenir la charge?



## Temps de réponse

### Percentile

## Utilisation mémoire

### Pression sur le Garbage Collector

## Utilisation CPU

Si une machine tombe, est-ce que les autres ont l'espace nécessaire pour soutenir la charge?



## Temps de réponse

### Percentile

## Utilisation mémoire

### Pression sur le Garbage Collector

## Utilisation CPU

Si une machine tombe, est-ce que les autres ont l'espace nécessaire pour soutenir la charge?



## Temps de réponse

### Percentile

## Utilisation mémoire

### Pression sur le Garbage Collector

## Utilisation CPU

Si une machine tombe, est-ce que les autres ont l'espace nécessaire pour soutenir la charge?

# Trucs intéressants



**99.999%**

# Les autres “9”s

“your nines are not my nines”



0.000%



99.999%

# Les autres “9”s

“your nines are not my nines”

QPS ▾



# Les autres “9”s

“your nines are not my nines”

**QPS**  
Query per seconds

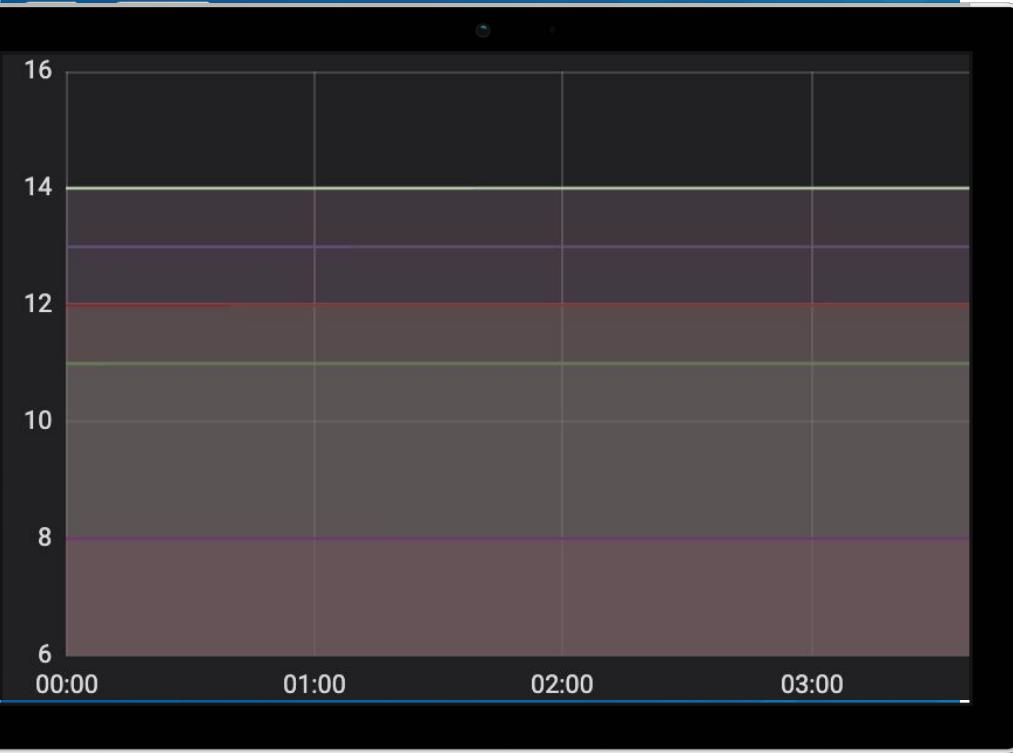
# **Les autres “9”s**

“your nines are not my nines”

**QPS**  
Query per seconds

**Invariant**  
observations extérieures

# Une histoire de **programmeur**

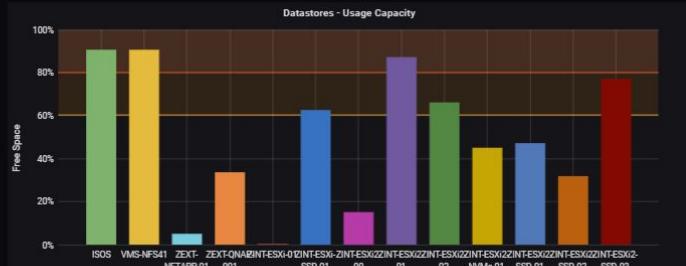
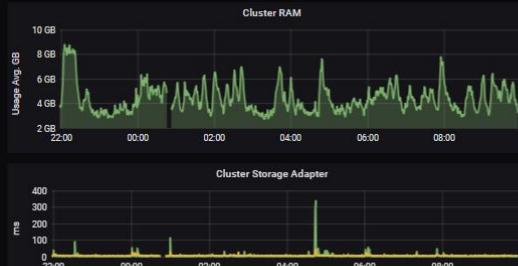
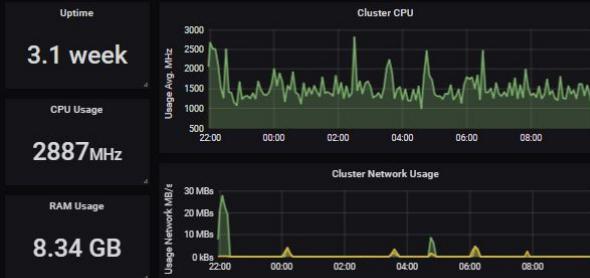


# Invariant

Surveiller un système complexe

# Image de l'état

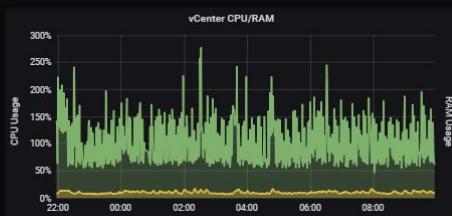
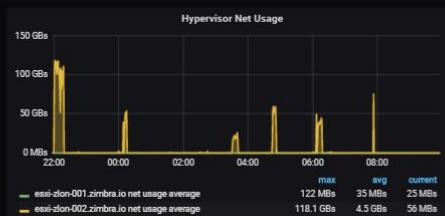
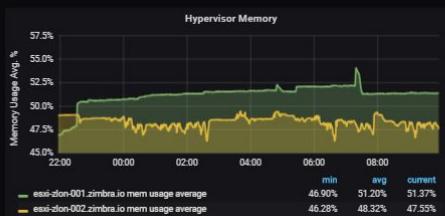
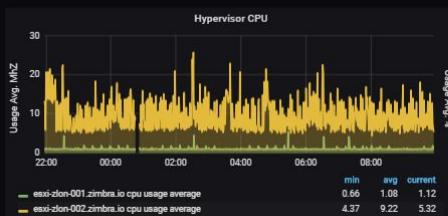
## Cluster Status



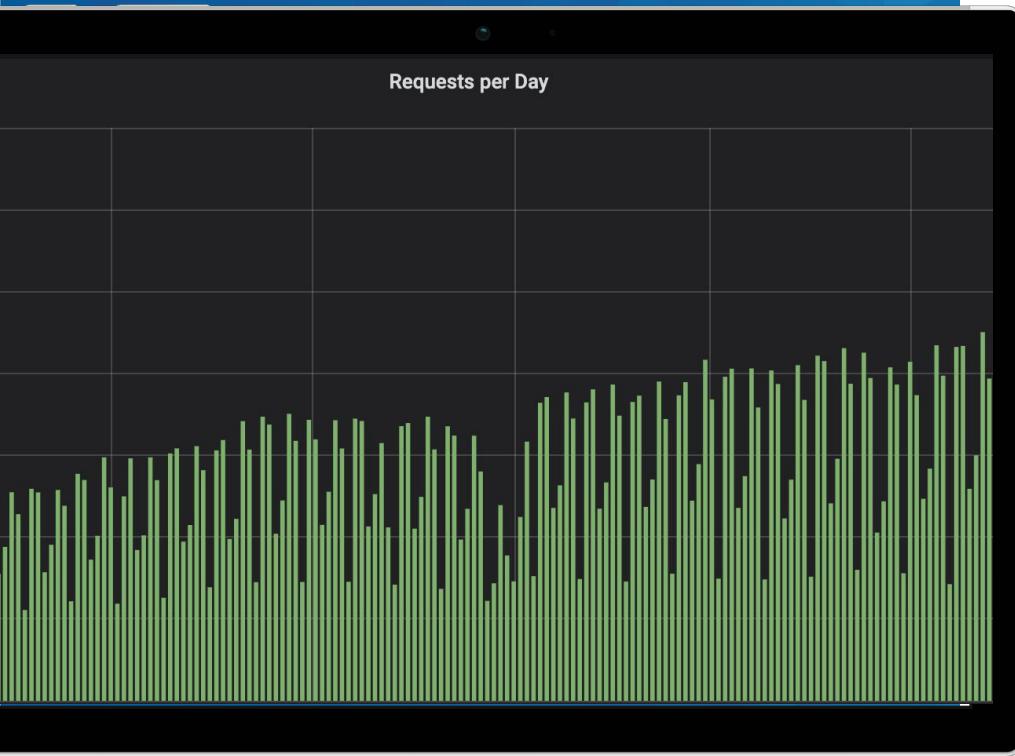
## Datastore Status



## Hypervisor Status



10  
QPS



# La situation.

13  
QPS



YEAH!!

Chaque  
semaine est la  
**plus occupée**

15  
QPS

Jusqu'à **doubler**  
**de volume** aux  
deux mois



30  
QPS

Quoi faire si nos métriques nous indiquent une faiblesse dans le système?

32  
QPS

# Nouveau déploiement!

35  
QPS

# Nouveau déploiement?

38  
QPS

# On déploie aux deux mois?

40  
QPS

On déploie aux deux mois?

Les soirs & fins de semaine?

43  
QPS

# On déploie aux deux mois?

Les soirs & fins de semaine?

Avec downtime pour les clients ??

A red toy car is shown from a top-down perspective, driving down a steep, rocky hillside. The hillside is covered in brown and tan rocks, and the ground is a mix of dirt and small stones. The car is positioned in the lower right quadrant of the frame, moving towards the bottom right.

46  
QPS

47

QPS

# Raccourcir la boucle de rétroaction

48  
QPS

# Raccourcir la boucle de rétroaction

## Pour les développeurs

- ▶ Cycle de développement plus court
- ▶ Master doit **toujours** être fonctionnel
- ▶ Tests plus rapide
- ▶ Compilation plus rapide
- ▶ Déploiement automatisé

# Raccourcir la boucle de rétroaction

## Pour les développeurs

- ▶ Cycle de développement plus court
- ▶ Master doit **toujours** être fonctionnel
- ▶ Tests plus rapide
- ▶ Compilation plus rapide
- ▶ Déploiement automatisé

## Pour les clients

52  
QPS

54  
QPS

# Raccourcir la boucle de rétroaction

## Pour les développeurs

- ▶ Cycle de développement plus court
- ▶ Master doit **toujours** être fonctionnel
- ▶ Tests plus rapide
- ▶ Compilation plus rapide
- ▶ Déploiement automatisé

## Pour les clients Wô.

55  
QPS

Instabilité est introduite  
par le **changement**

57  
QPS

Mean time  
to failure

vs

Mean time  
to recovery

**58**  
QPS

Mean time  
to failure

VS

Mean time  
to recovery

- “Les pannes sont évitables”
- Maximise le temps entre les échecs
- Veut attraper tous les bugs en amont (long QA, processus, etc)
- Veut minimiser le risque lors de moments critiques (brownout, blackout)

**60**  
QPS

## Mean time to failure

**VS**

## Mean time to recovery

- “Les pannes sont évitables”
- Maximise le temps entre les échecs
- Veut attraper tous les bugs en amont (long QA, processus, etc)
- Veut minimiser le risque lors de moments critiques (brownout, blackout)

- “Les pannes existent, réduisons les impacts de celles-ci”
- Minimise le temps des pannes
- Minimise l’impact des pannes

# Raccourcir la boucle de rétroaction

## Pour les développeurs

- ▶ Cycle de développement plus court
- ▶ Master doit **toujours** être fonctionnel
- ▶ Tests plus rapide
- ▶ Compilation plus rapide
- ▶ Déploiement automatisé

63  
QPS

# Raccourcir la boucle de rétroaction

## Pour les développeurs

- ▶ Cycle de développement plus court
- ▶ Master doit **toujours** être fonctionnel
- ▶ Tests plus rapide
- ▶ Compilation plus rapide
- ▶ Déploiement automatisé

## Pour les clients

- ▶ Automatisation des déploiements

64  
QPS

# Raccourcir la boucle de rétroaction

66  
QPS

## Pour les développeurs

- ▶ Cycle de développement plus court
- ▶ Master doit **toujours** être fonctionnel
- ▶ Tests plus rapide
- ▶ Compilation plus rapide
- ▶ Déploiement automatisé

## Pour les clients

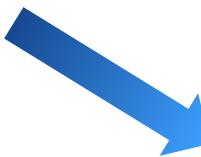
- ▶ Automatisation des déploiements
- ▶ Cacher les nouvelles fonctionnalités derrière *feature flags*

70  
QPS

Old

New

```
if (enabled)
```



# Raccourcir la boucle de rétroaction

72  
QPS

## Pour les développeurs

- ▶ Cycle de développement plus court
- ▶ Master doit **toujours** être fonctionnel
- ▶ Tests plus rapide
- ▶ Compilation plus rapide
- ▶ Déploiement automatisé

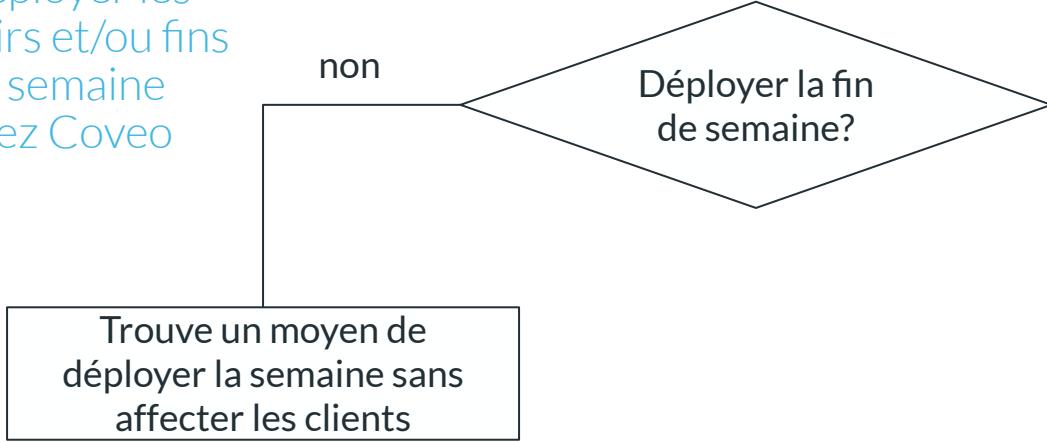
## Pour les clients

- ▶ Automatisation des déploiements
- ▶ Cacher les nouvelles fonctionnalités derrière *feature flags*
- ▶ *Canary releases*

74  
QPS



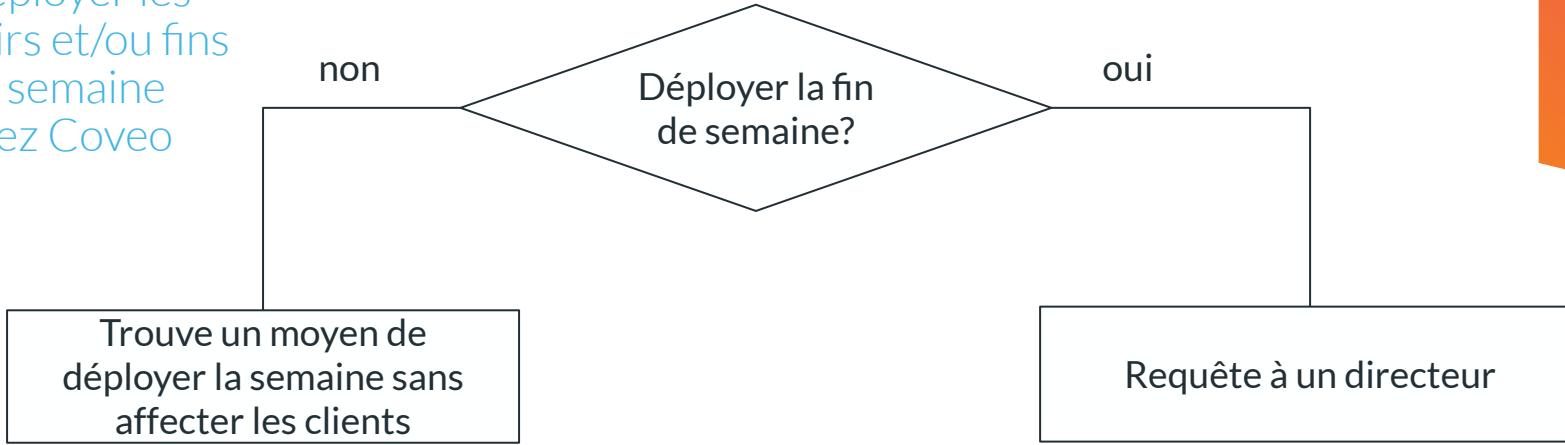
## Déployer les soirs et/ou fins de semaine chez Coveo



77  
QPS

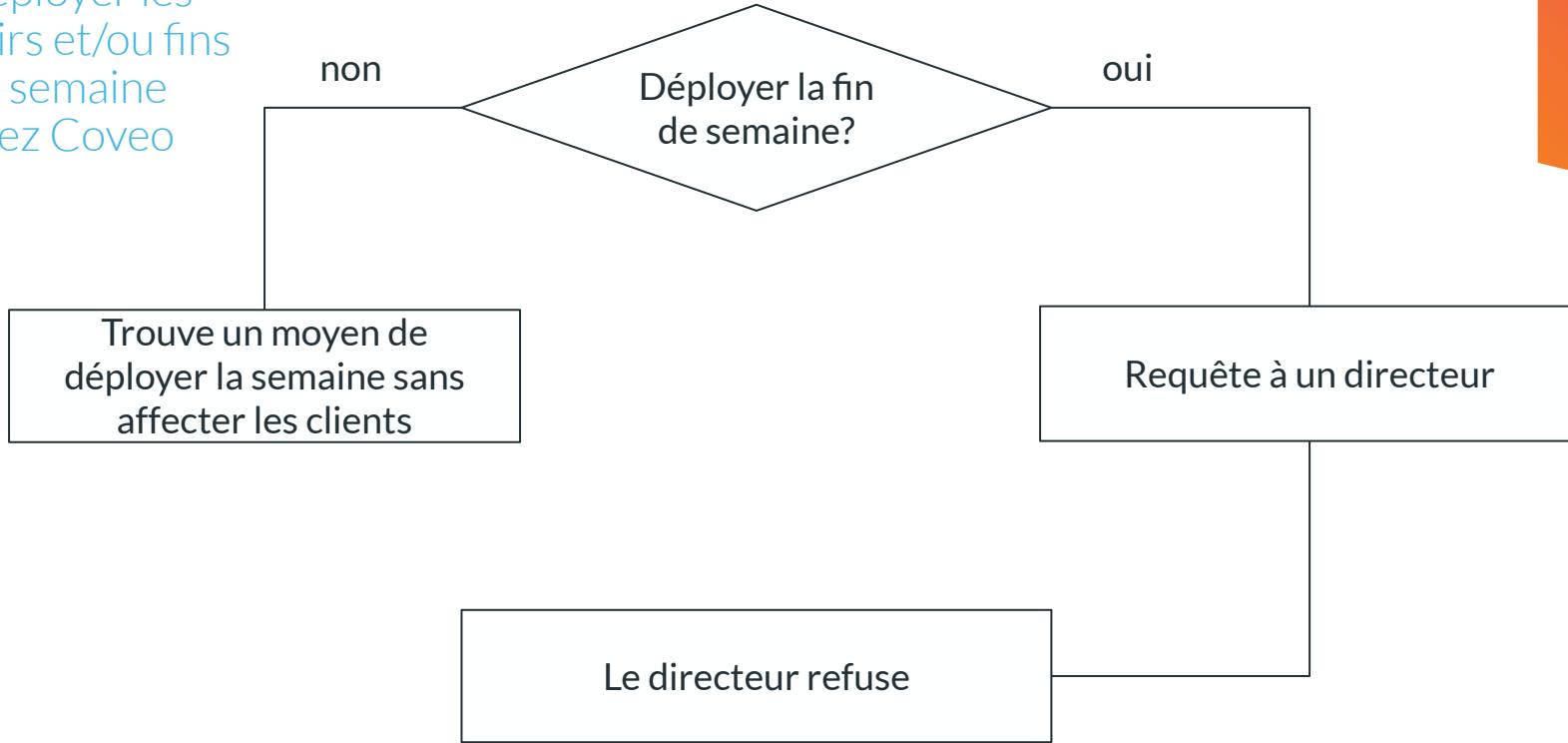
## Déployer les soirs et/ou fins de semaine chez Coveo

78  
QPS



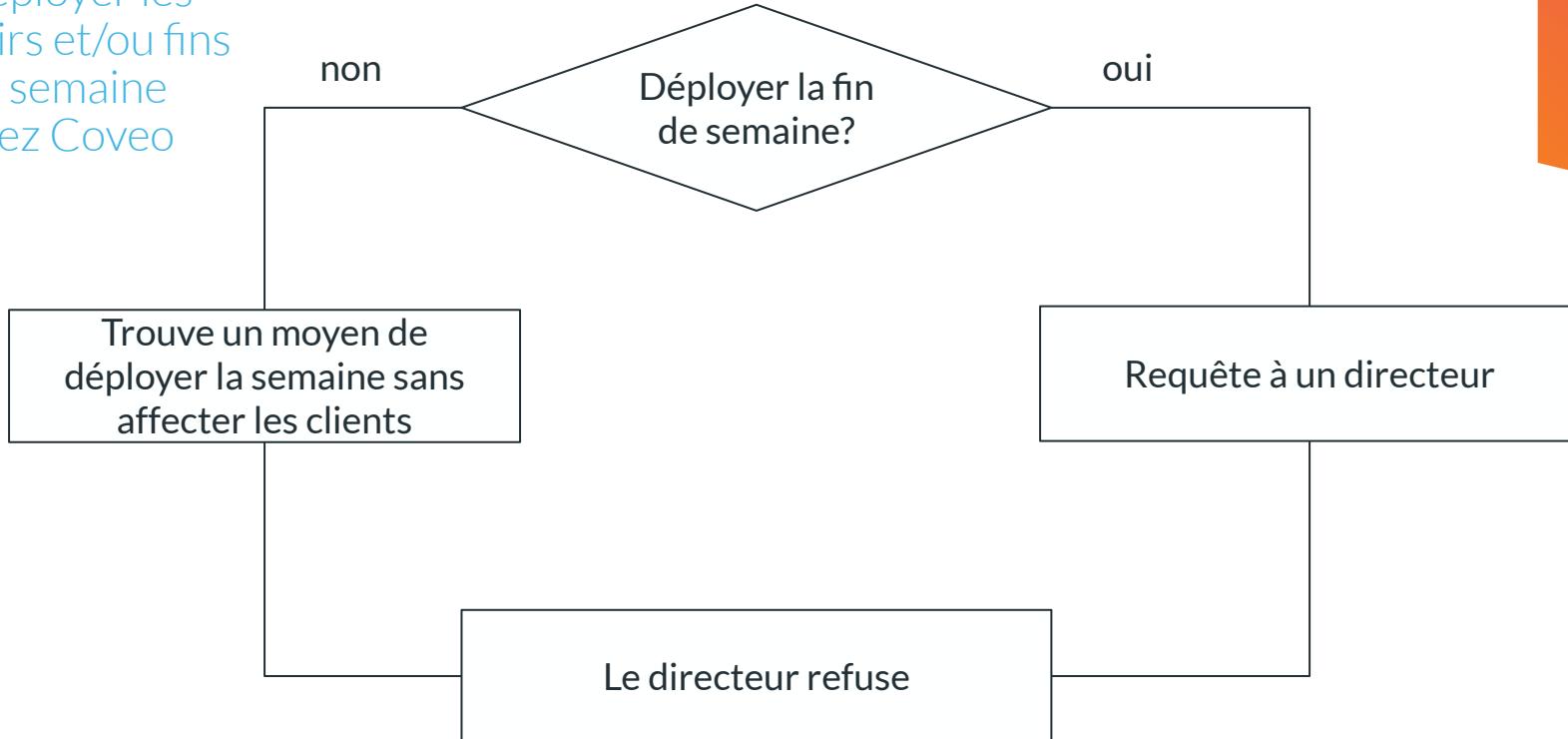
# Déployer les soirs et/ou fins de semaine chez Coveo

82  
QPS

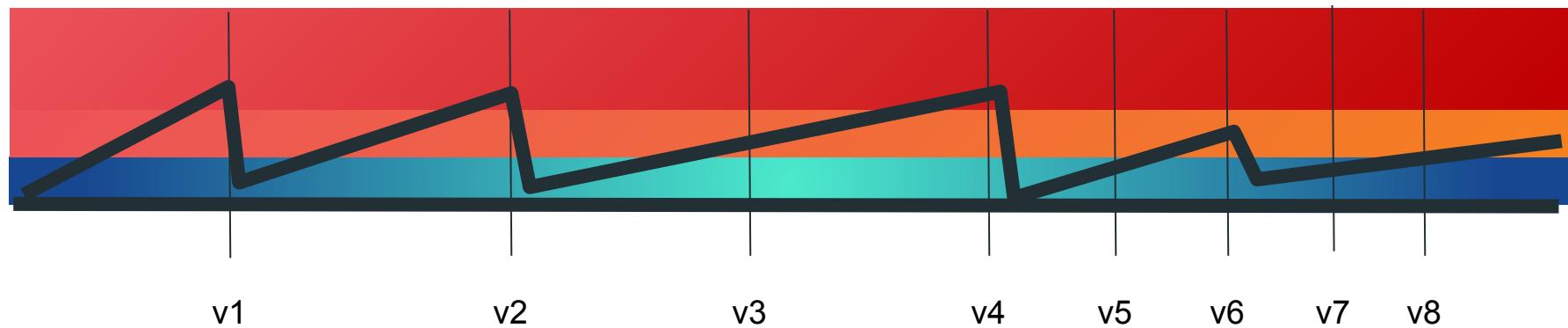


# Déployer les soirs et/ou fins de semaine chez Coveo

85  
QPS



90  
QPS



96  
QPS

# Comment développer de grandes fonctionnalités?

# Comment développer de grandes fonctionnalités?

## Ré-écriture

- ▶ Pas une option pour nous
- ▶ Historiquement plus cher, plus long, plus compliqué que prévu (Netscape, Parse, etc)

99

QPS

# Comment développer de grandes fonctionnalités?

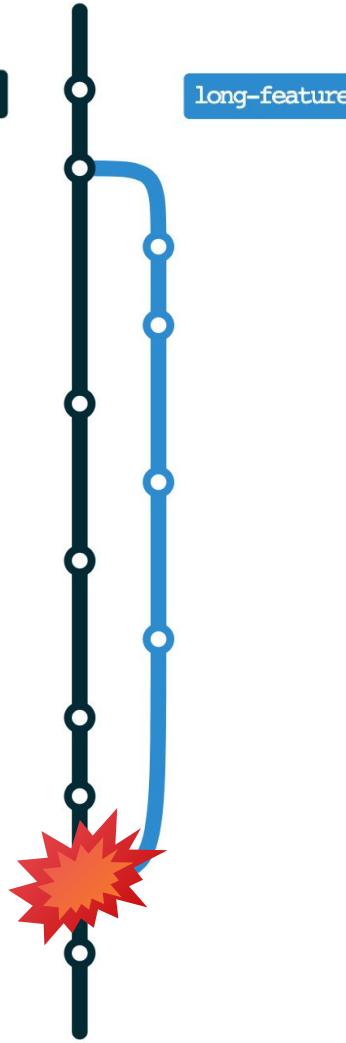
## Ré-écriture

- ▶ Pas une option pour nous
- ▶ Historiquement plus cher, plus long, plus compliqué que prévu (Netscape, Parse, etc)

## Long lived branch

- ▶ On a déjà mentionné que si c'est pas testé en production, ce n'est pas vraiment testé

103  
QPS



# Comment développer de grandes fonctionnalités?

## Ré-écriture

- ▶ Pas une option pour nous
- ▶ Historiquement plus cher, plus long, plus compliqué que prévu (Netscape, Parse, etc)

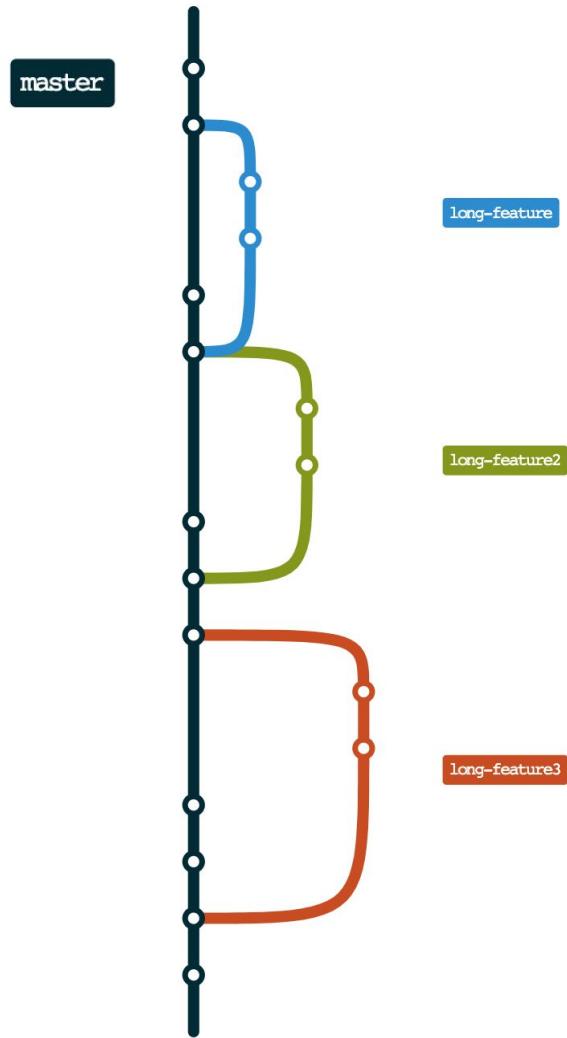
## Long lived branch

- ▶ On a déjà mentionné que si c'est pas testé en production, ce n'est pas vraiment testé

## Morcelé + caché

- ▶ *Feature Flag*
- ▶ *Phased deployment*

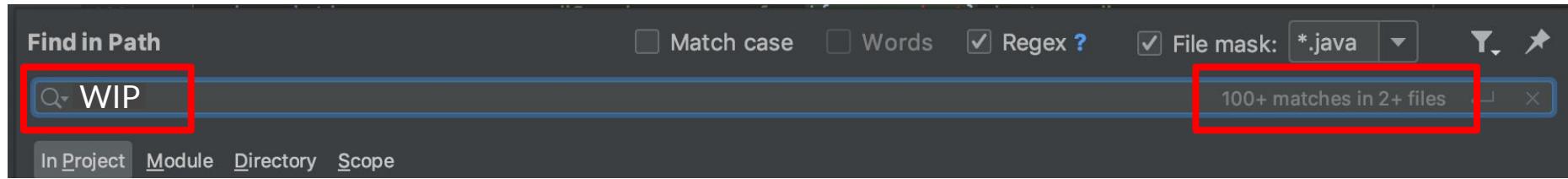
**115**  
QPS



125  
QPS

Déploiement en phases  
+ changement minimal ?

130  
QPS

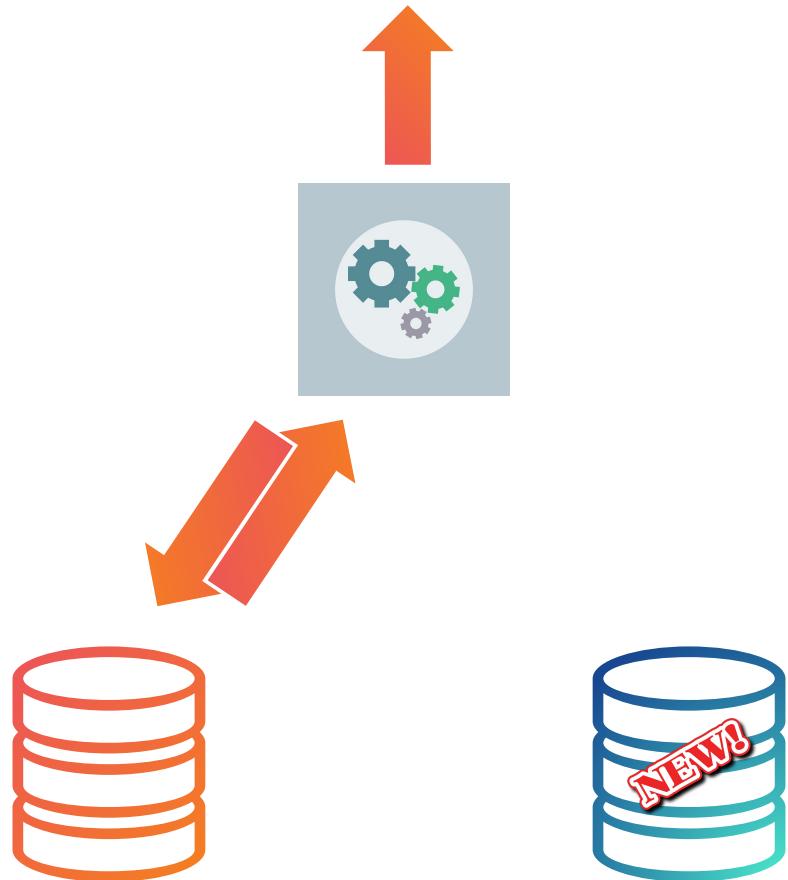


135  
QPS

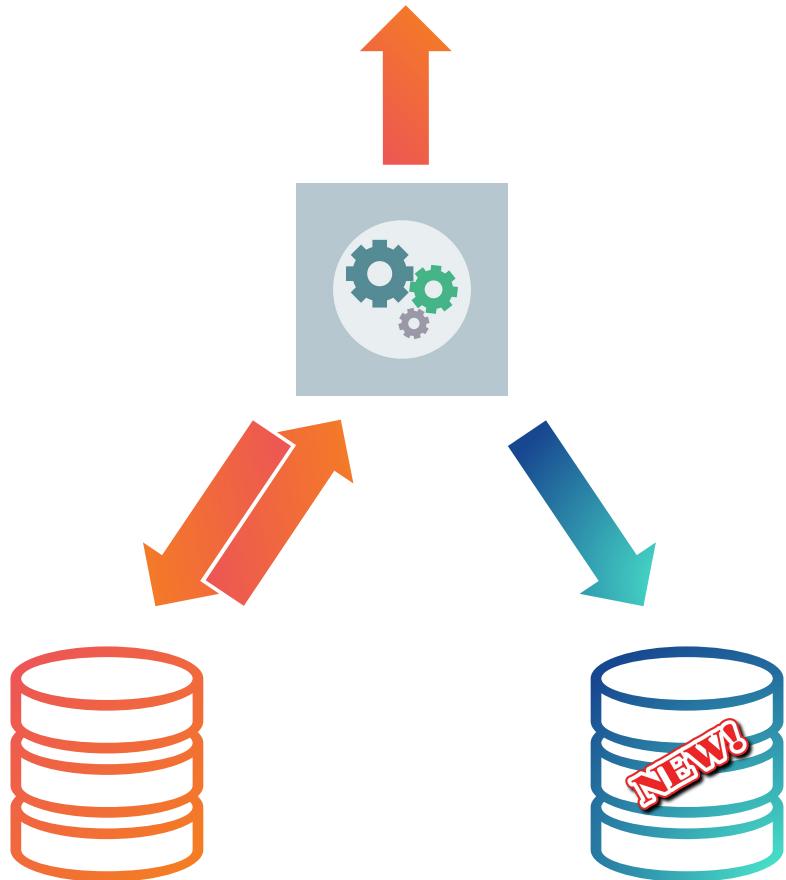
## Exemple

# Migrer des données vers une nouvelle base de données

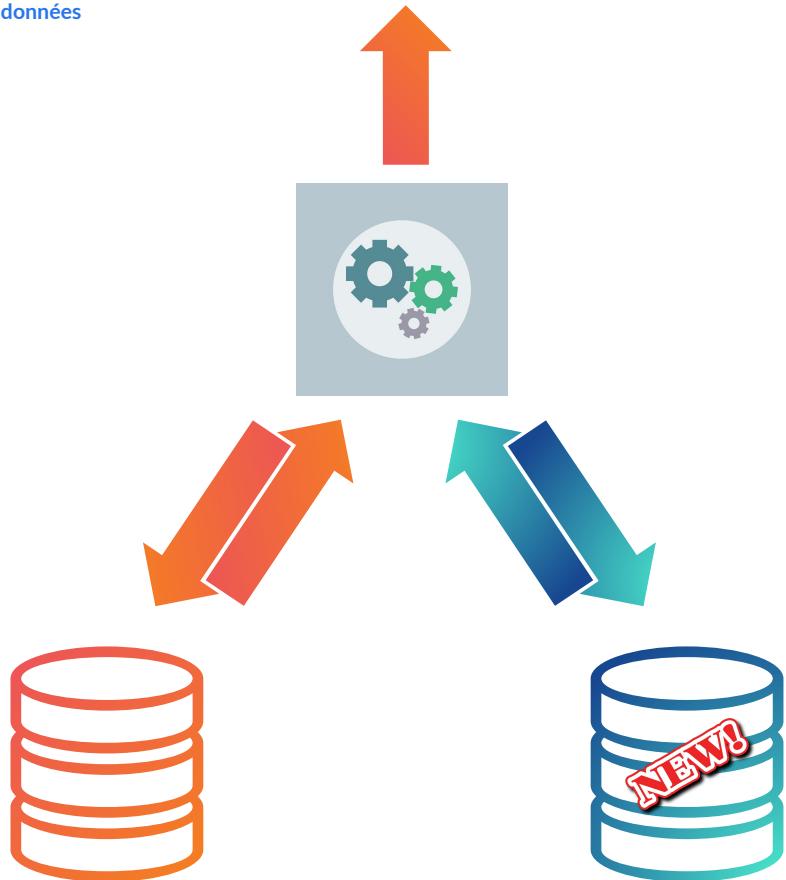
137  
QPS



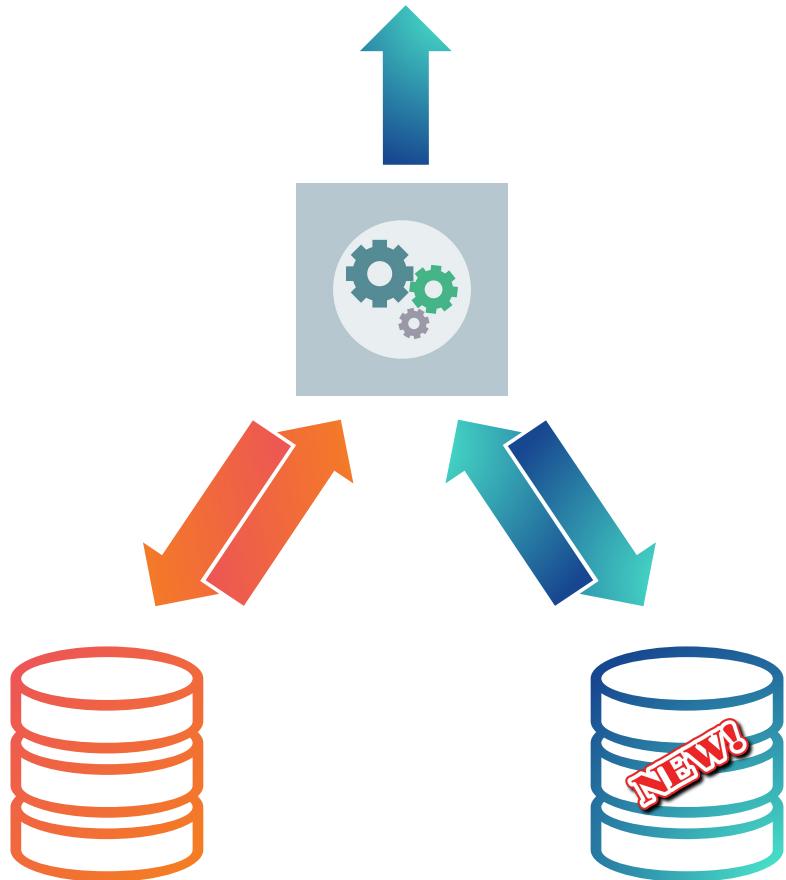
138  
QPS



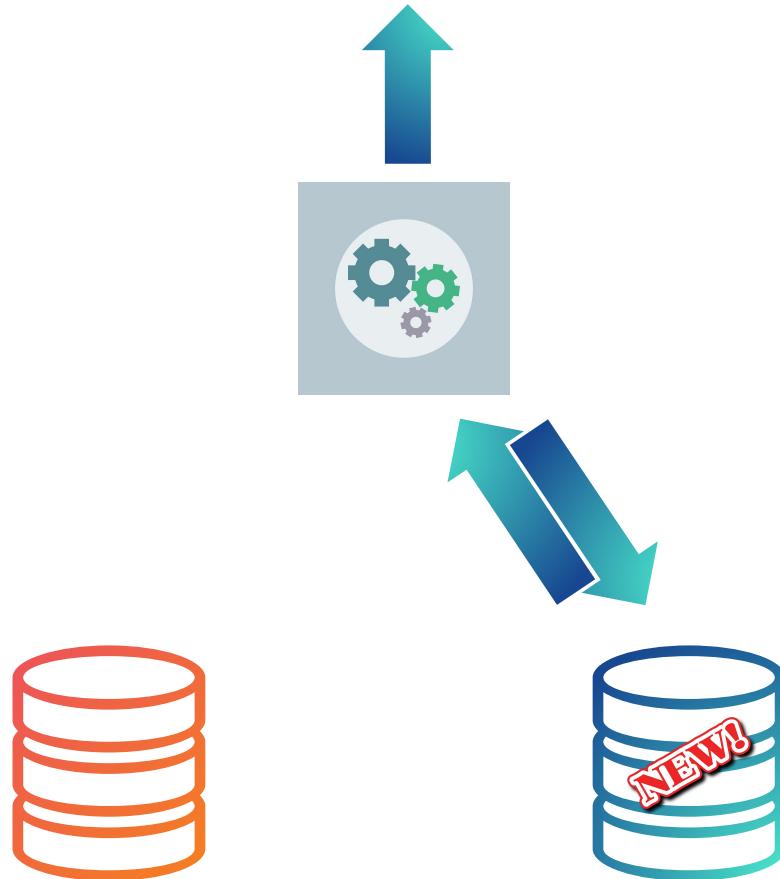
135  
QPS



141  
QPS



144  
QPS



# (Lightweight) Architecture Decision Record

## ADR

- ▶ Simple
- ▶ Avec le code
- ▶ Versionné
- ▶ *git blame friendly <3*

```
1  # Some Performance issue about something #
2  Architecture decision record for performance issue
3
4  ### Status ###
5  Accepted in [MyProject-1759](https://myproj.atlassian.net/b...
6
7  ### Context ###
8
9  * We are getting throttled
10 * We are using a O(n^2) algorithm
11 * Some of our cached methods called other cached methods, b...
12
13 ### Decisions ###
14
15 * Temporary use a local cache
16 * Add different caching durations depending on the informat...
17 * Include an abstraction layer for all EMR calls
18
19 ### Diagram ###
20 [Beautiful Diagram](../diagrams/beautiful.xml)
21
```

160  
QPS

# Comment **planifier** le travail?

**164**  
QPS

BEST BEFORE  
**300 QPS**

167  
QPS

# Metric Driven Development

168

QPS

# Zinga



“La croissance, ça ne se règle pas en ajoutant plus de monde, mais en priorisant bien les tâches.”

MAXIME LACHAPELLE

177  
QPS

Investir sur  
le vieux      **vs**

Attendre  
le neuf



QPS

Conclusion ou quelque chose de  
pertinent du genre



# Thank you



