

A decorative graphic consisting of several parallel diagonal lines in a light green color, slanted from the top-left to the bottom-right.

Construindo Uma Plataforma Data-intensive do Zero

Um estudo de caso na In Loco

Vinícius M.R. Cousseau

Engenheiro de software na In Loco, trabalhando há 3 anos no time de R&D em localização. Mestrando no CIn-UFPE, pesquisando na área de deep learning aplicada a record linkage. Trabalhou previamente com Subsurface Scattering.

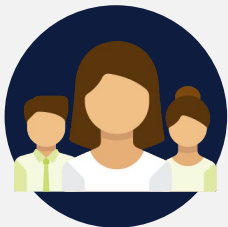
 [viniciuscousseau](#)



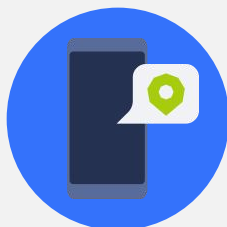
Objetivo

Mostrar um roadmap da nossa experiência na In Loco ao desenvolver uma plataforma data-intensive do zero, para gerar insights e direcionamentos.

Quem somos?



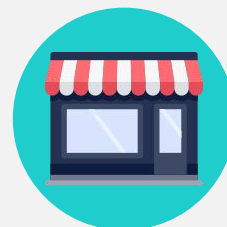
60 MM
usuários únicos
em nossa base



16 TB
de dados colhidos
diariamente



1.5 BI
visitas
registradas
mensalmente



+28MM
locais
mapeados no
mundo



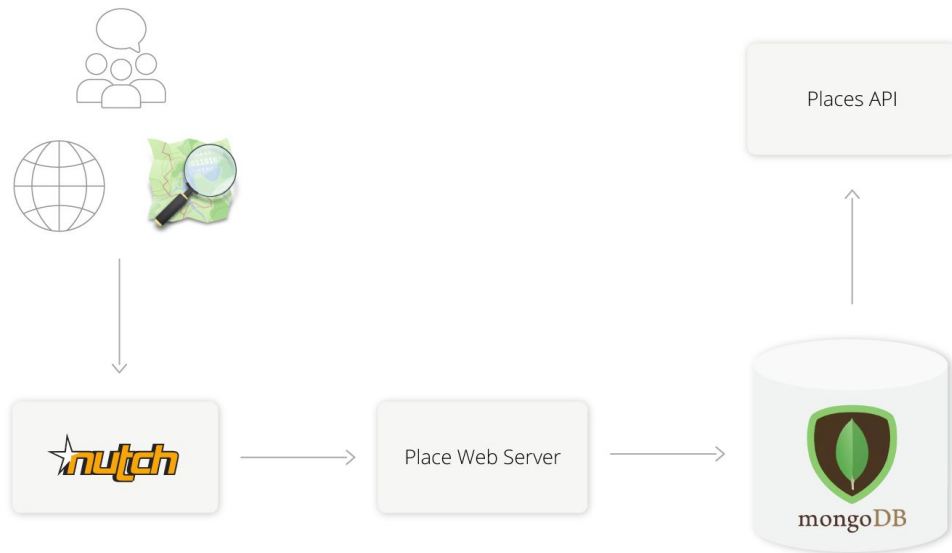
Desafios

- Mapear os locais do mundo
- Fornecer esses dados através de uma plataforma
- Processamento On-line + Off-line
- Third party é inviável



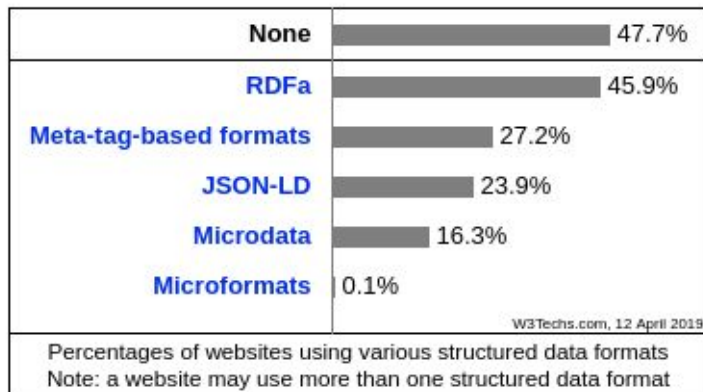
First Steps

Data Extraction



Data Extraction

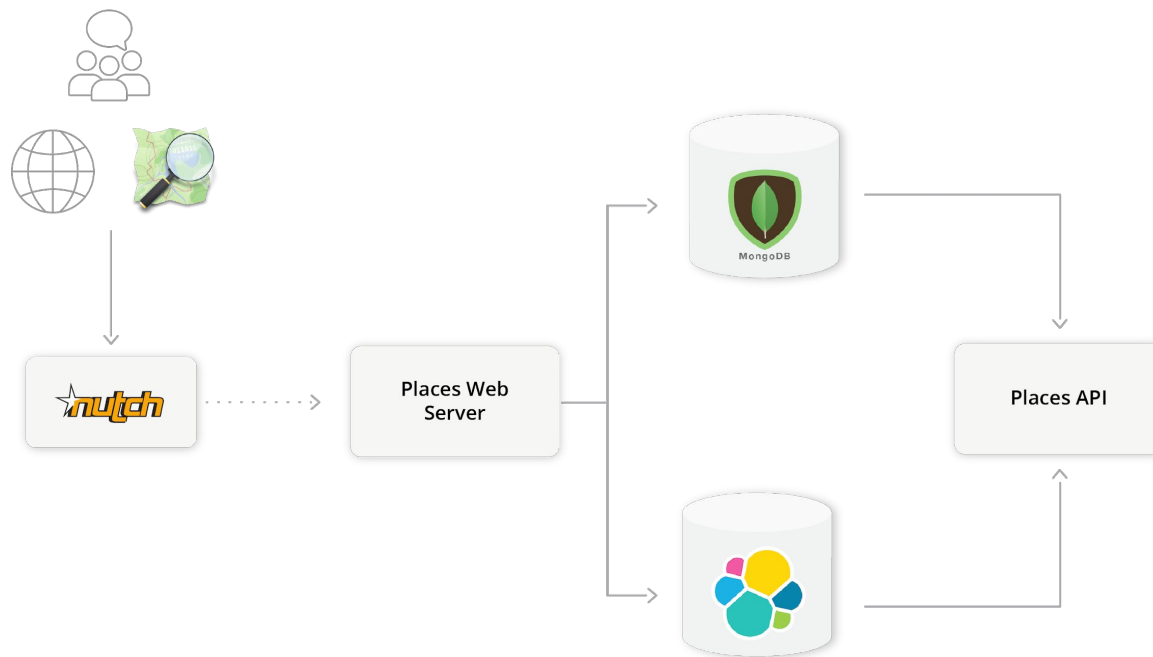
- ***Focused Crawling* com *nutch* plugins de dados estruturados (sempre respeitando as regras)**



Data Extraction

- 120k records por dia
- 5k RPM na Places API
- Apenas duas rotas
- Expor desde o t0 nos trouxe muitos ganhos

Escolha de Bancos

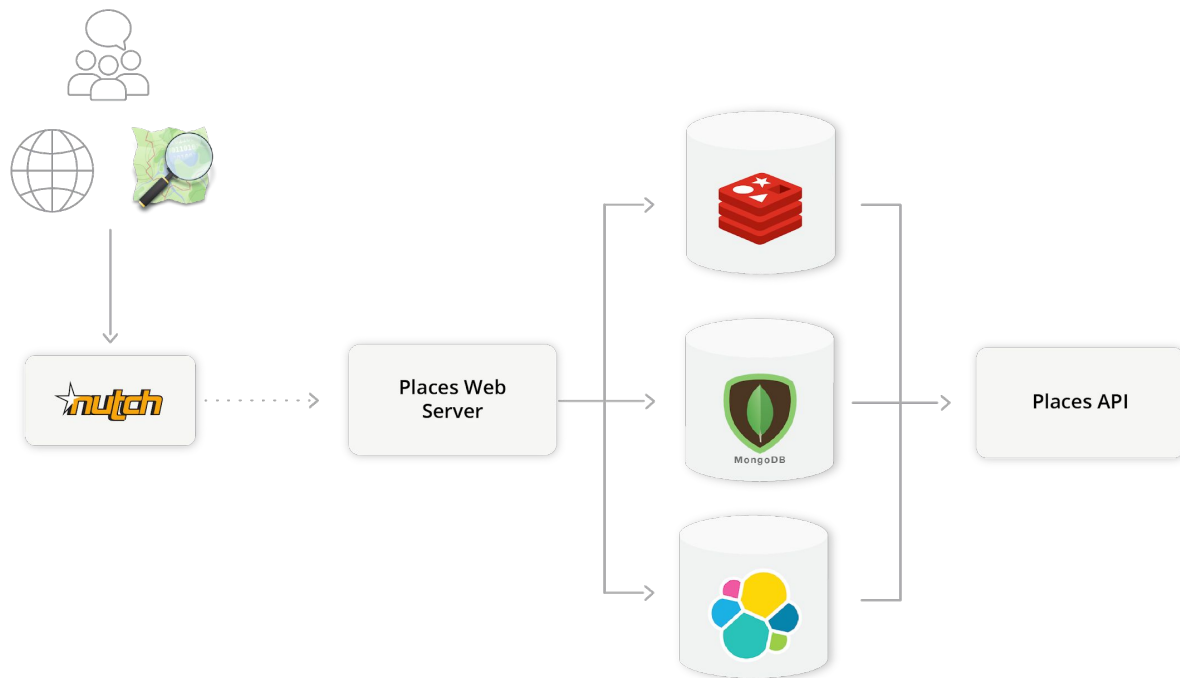


Escolha de Bancos

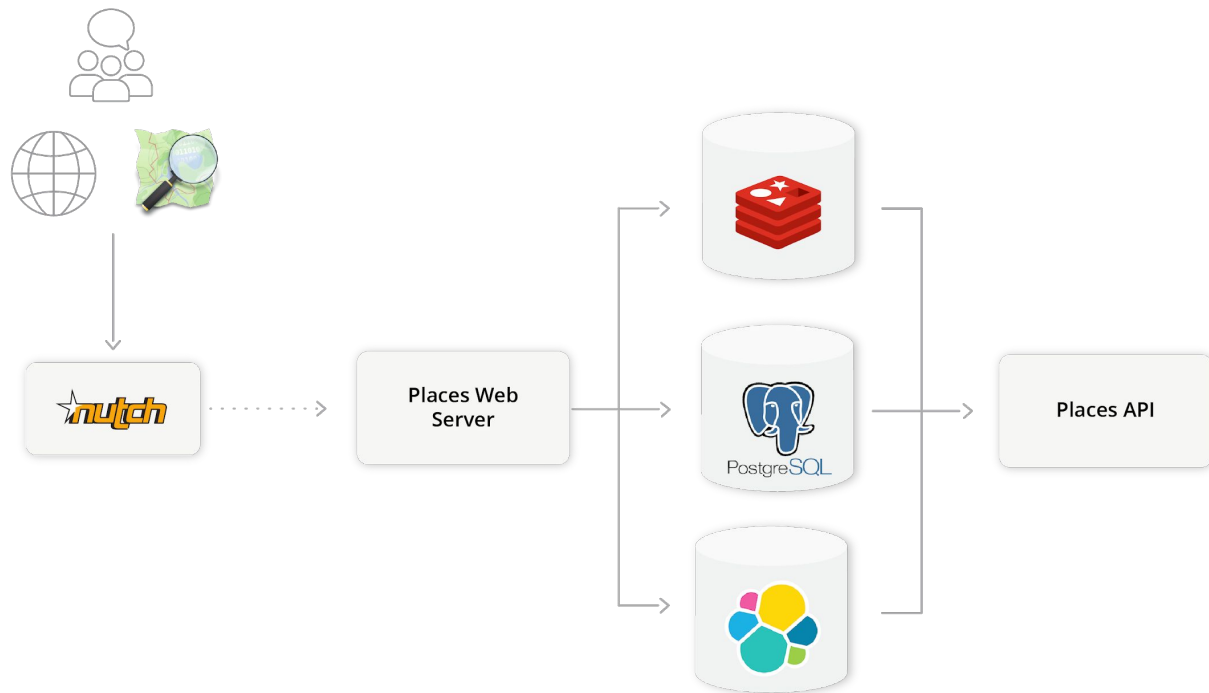
Database	Mean Memory	Mean CPU	Mean Throughput
Redis 3.2.9	5.2% (420MB)	50%	44438.7 rps
Mongo 3.2.9	2% (164MB)	100%	5676.9 rps

Geo Queries dinâmicas em um raio de 200m e limite de 100 resultados. Banco com ~25m records

Escolha de Bancos

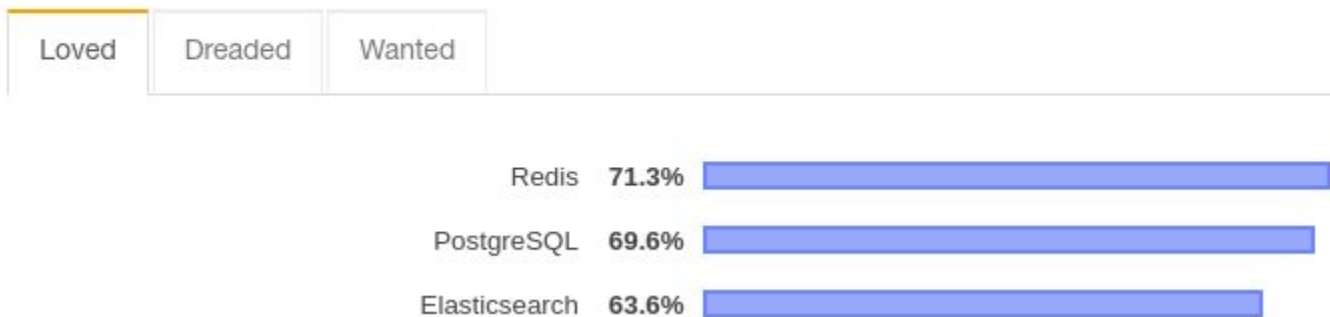


Escolha de Bancos



Escolha de Bancos

Most Loved, Dreaded, and Wanted Databases



Escolha de Bancos

- O domínio dos dados é um forte influenciador
- *No silver bullet* (out of the box)
- *Benchmark!*



Processamento em Batch

Processamento off-line

- *More domain knowledge, more problems*
- Web Server com dificuldades em picos

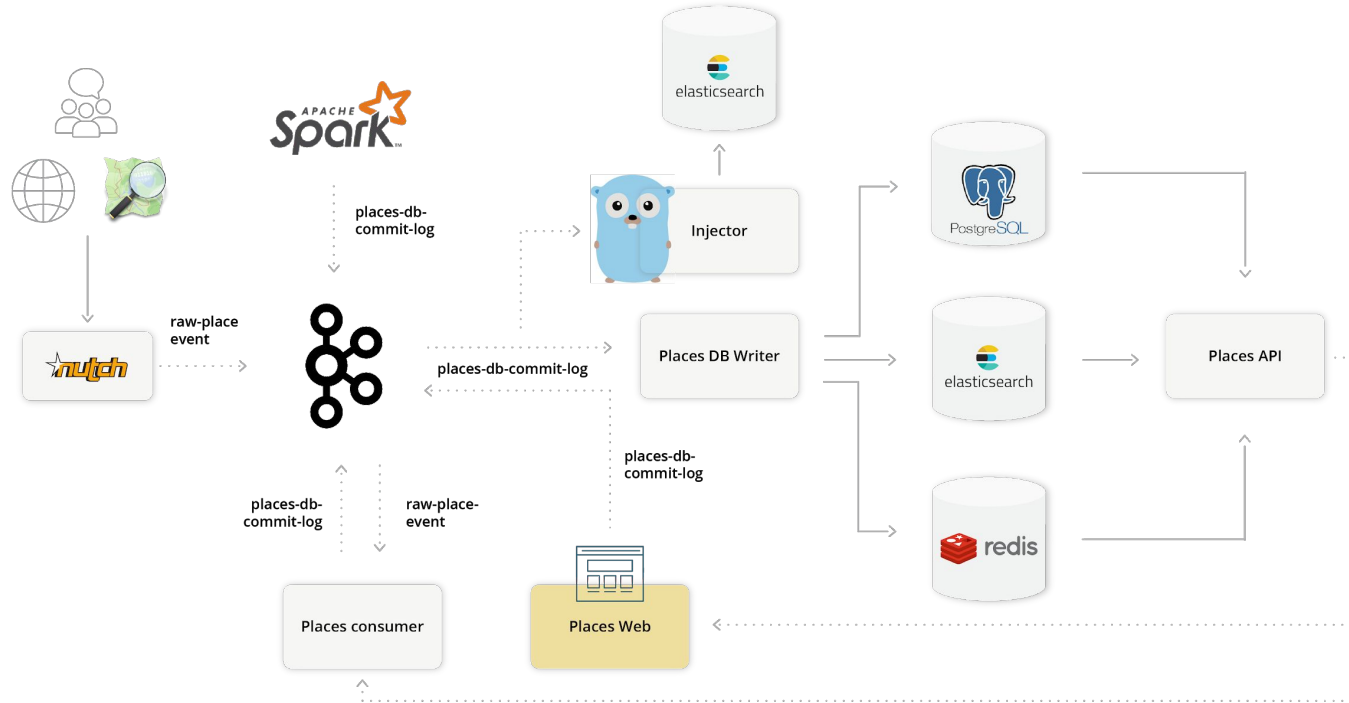
Processamento off-line



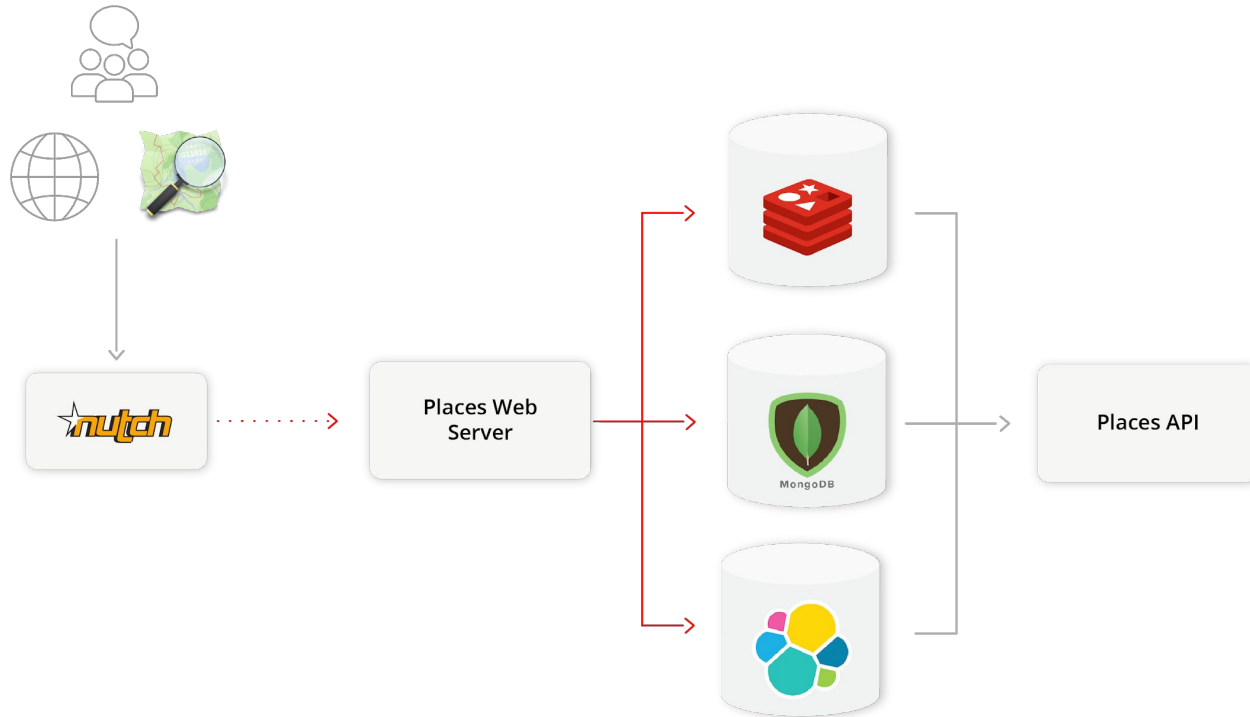
Apache Spark + Kafka

- **Retira nossos limitantes de tempo**
- **Alinhamento com o requisito inicial**
- **O quão antes for feita essa parte, mais cedo será possível iterar sobre os dados**

Versão Atual



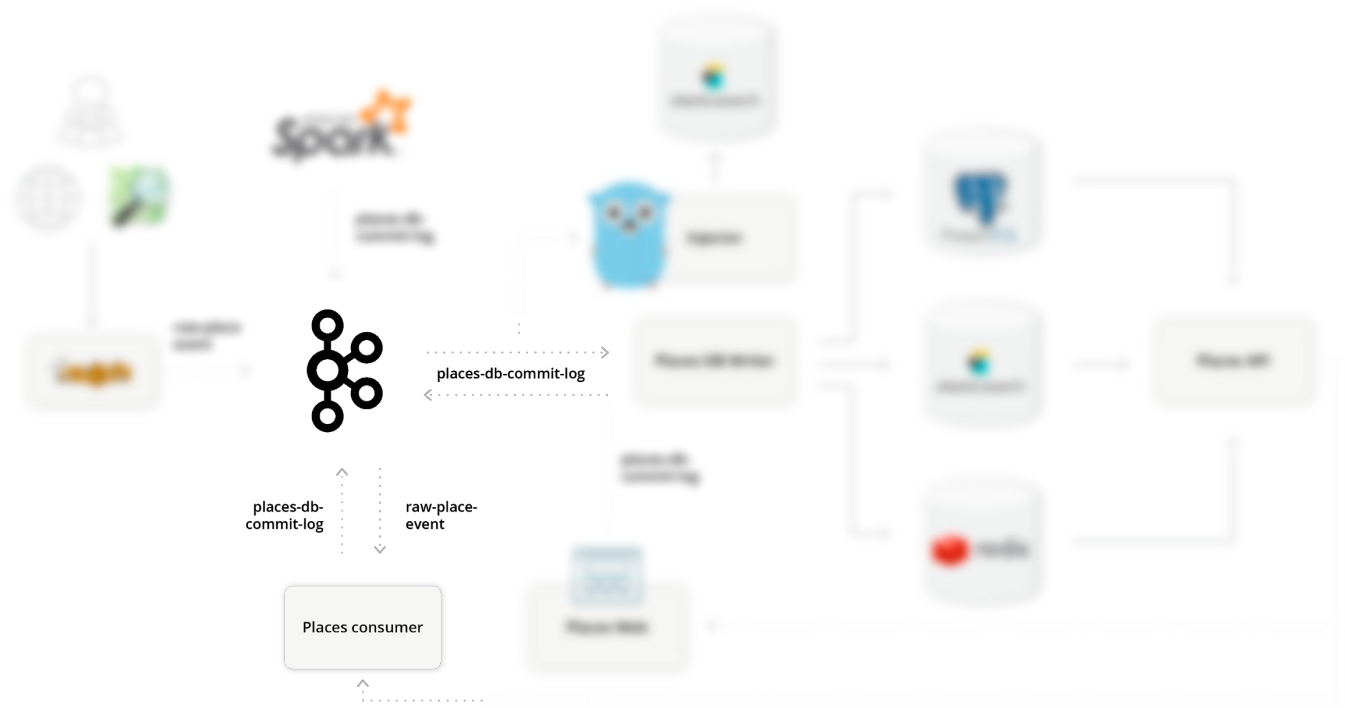
Remoção de Gargalos



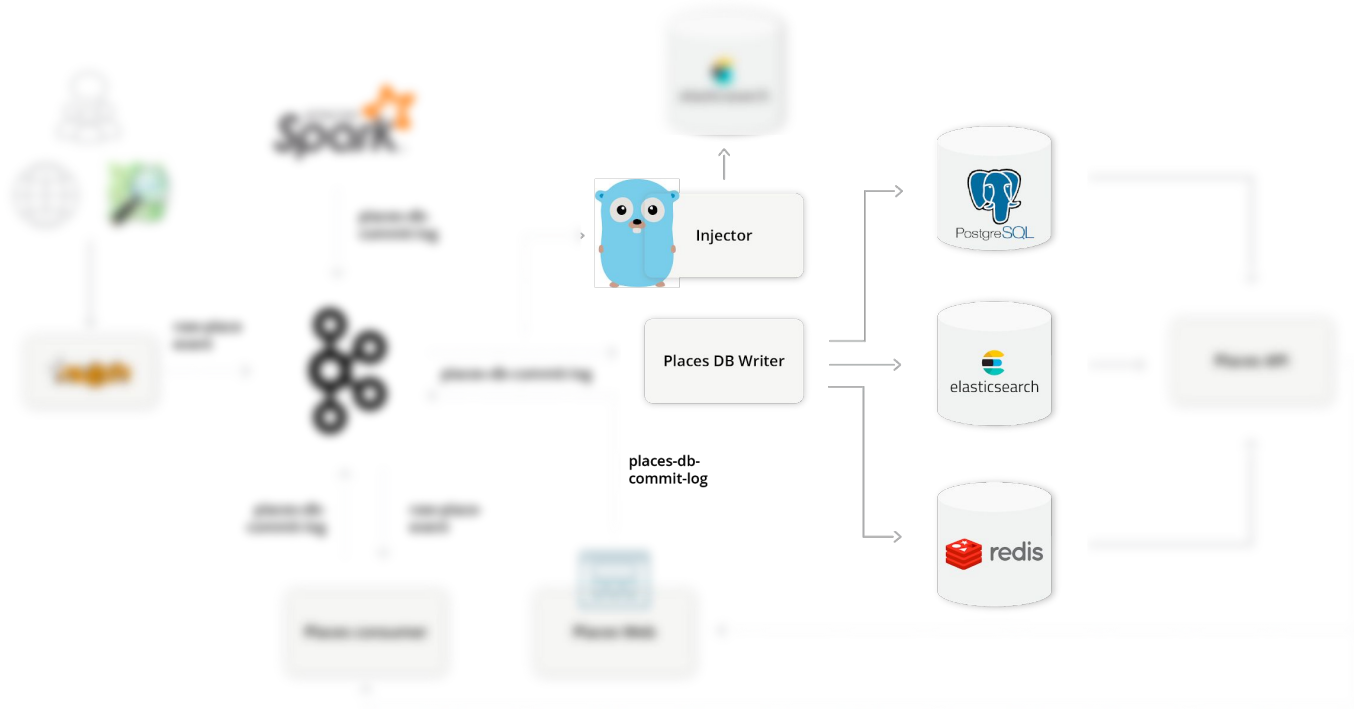
Remoção de Gargalos



Eventos



Eventos



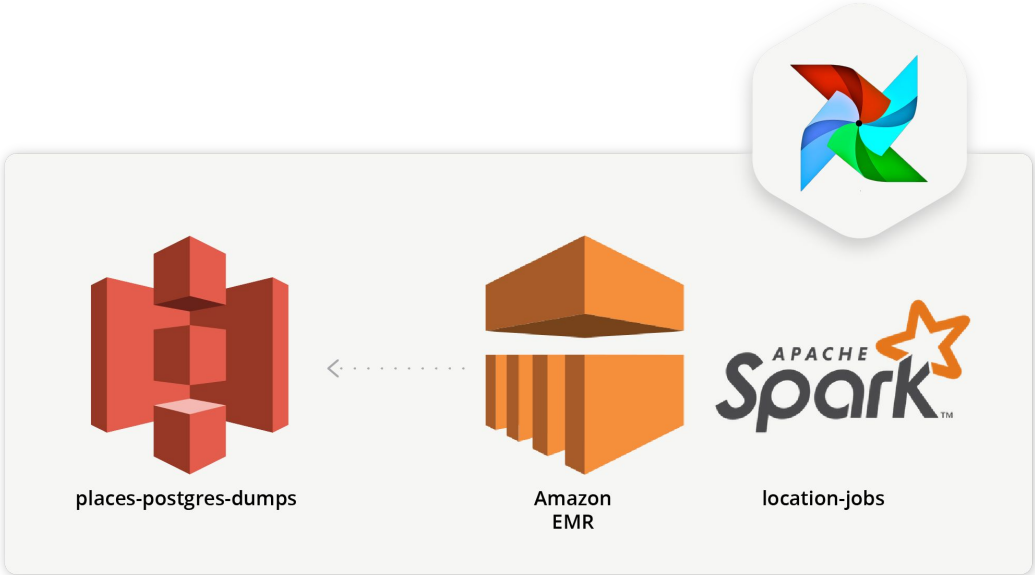
Eventos

- Replicação de CRUD nos 3 bancos
- Reversibilidade
- Visibilidade - Crucial!
 - *Can't improve what you cannot see*

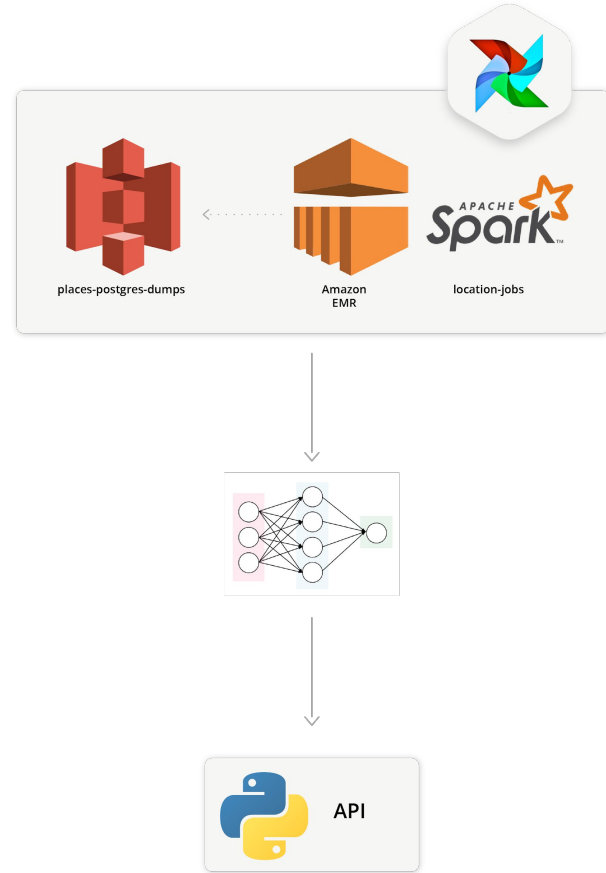
Visibilidade - Injector



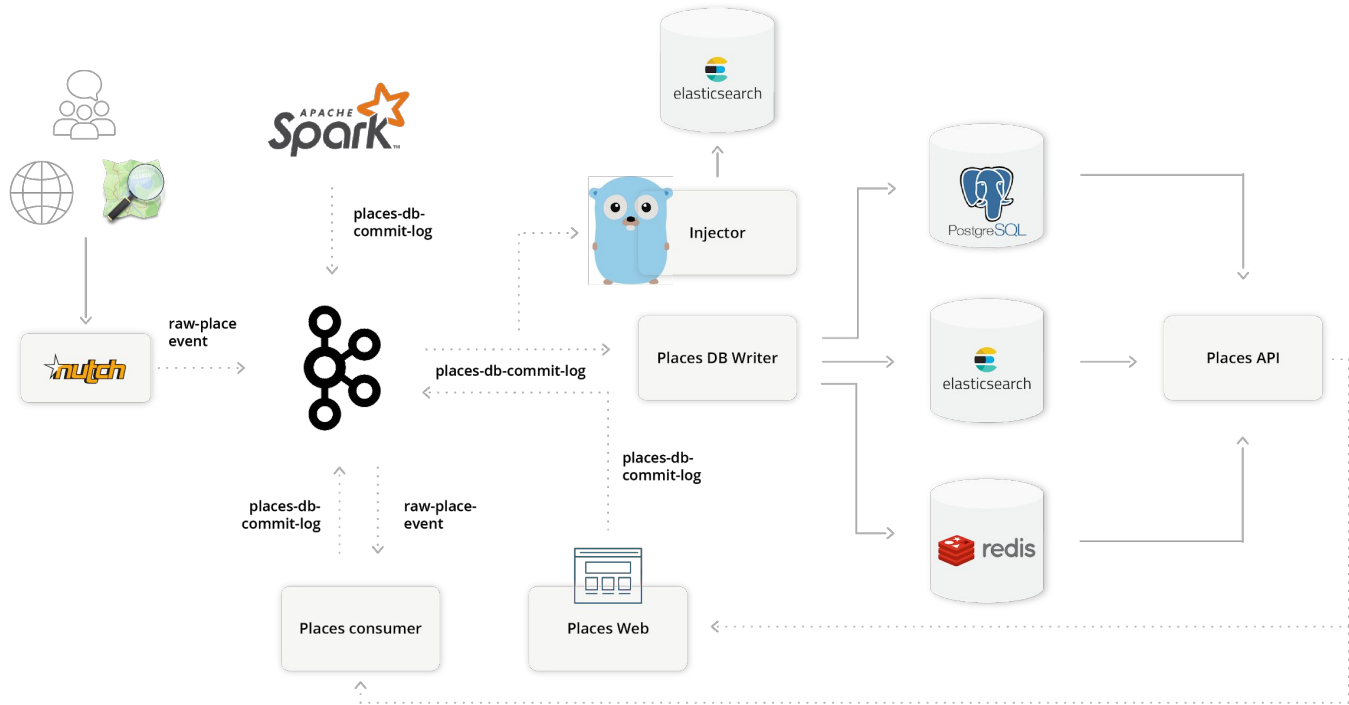
Orquestrando Jobs



Machine Learning



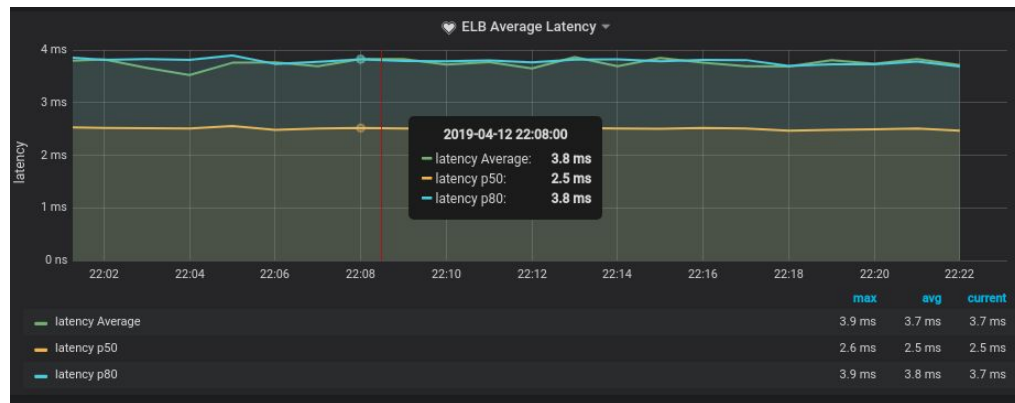
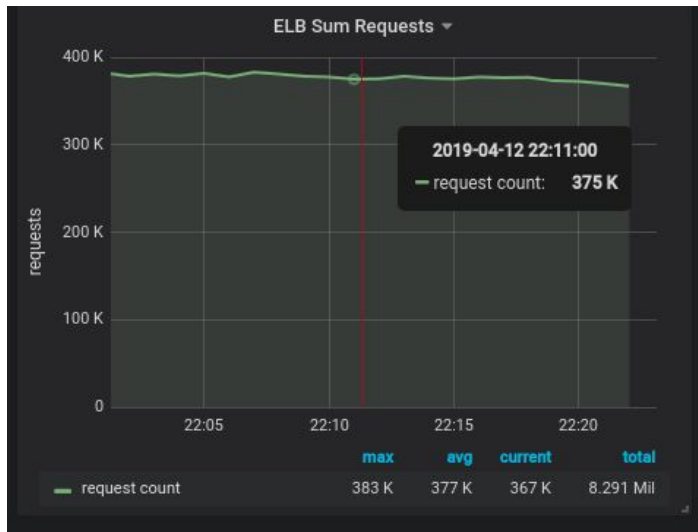
Versão Atual



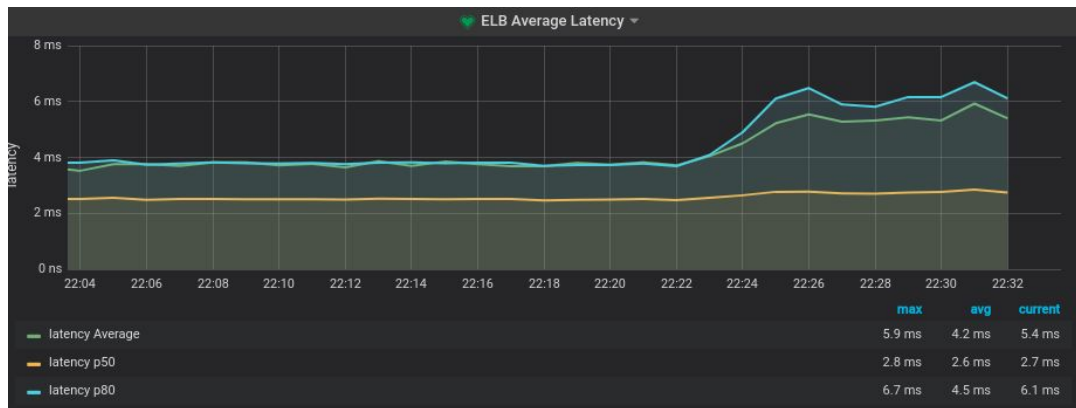
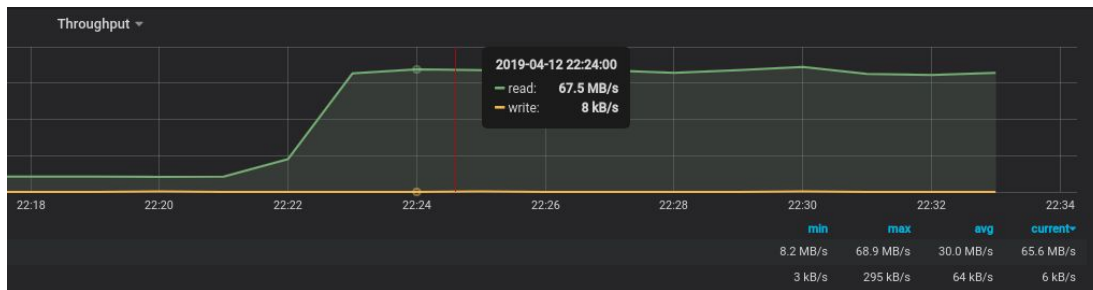


Outcomes

Escala



Escala



Sandbox

- 10 jobs rodando semanalmente
- Time de BI extraíndo insights diretamente dos nossos dados
- Em 1 mês, conseguimos iterar um dos nossos modelos (*record linkage*) 11 vezes



Onde Erramos

You can't connect the dots looking forward; you can only connect them looking backwards. So you have to trust that the dots will somehow connect in your future.

Steve Jobs



Pontos de Melhora

- Forçamos demais o uso de algoritmos distribuídos
- Fizemos as frentes on-line e off-line em silos
 - Atrasamos o benefício de *active learning*
- Demoramos a implementar o ferramental para reverter ações

Pontos de Melhora

- **Muito trabalho a se fazer ainda!**



vinicius.cousseau@inloco.com.br
inloco.com.br

Saiba mais em:
medium.com/inlocotech

Quer fazer parte do nosso time?
<https://inloco.com.br/en/careers>

Follow us

 [/inlocoglobal](#)

 [@inlocoglobal](#)

 [/inlocotech](#)