

Retour d'EXpérience ElasticStack

Jean-Pierre Paris
Orange
31 mai 2016



plan

- hier
- aujourd'hui
- demain

un moteur ou des moteurs

- moteurs thématiques
 - technologies variées
 - sous ES depuis janvier 2014
 - ça augmente !

- moteur web
 - technologie propriétaire
 - migration...

volumétrie : 1,2 mds urls
performances : 25 req/s, < 200

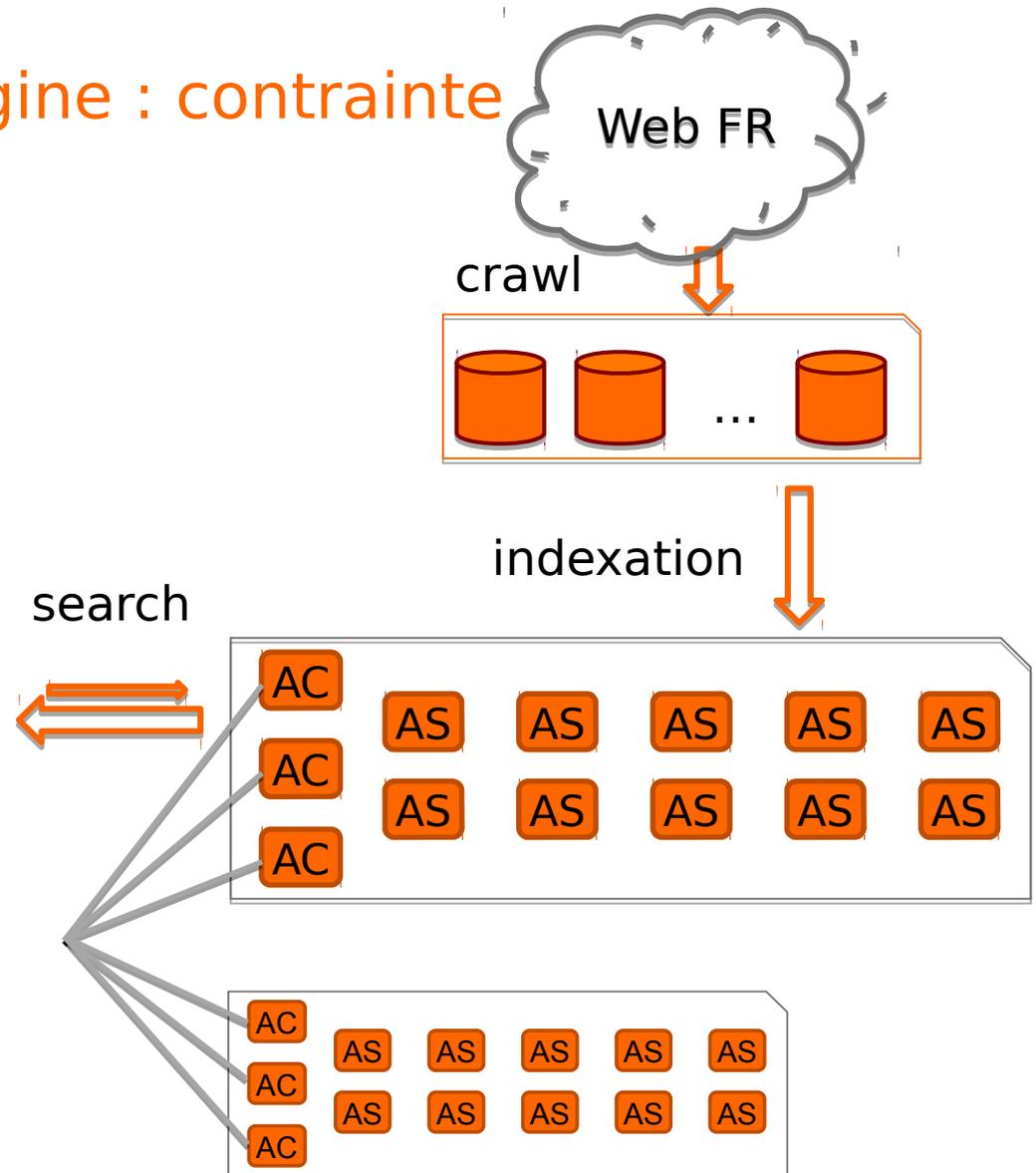
The screenshot shows the Orange search engine interface. At the top, there's a navigation bar with 'actualités', 'TV', 'VOD', 'jeux', and 'musique'. The search bar contains 'karim benzema'. Below the search bar, there are tabs for 'WEB', 'VIDÉOS', 'IMAGES', and 'ACTUALITÉS'. The main content area is divided into several sections:

- POUR ALLER VITE, C'EST LA**: A yellow box containing the 'Site officiel de Karim Benzema' with a small image and the URL 'benzema.com'.
- TOUTES LES NEWS ORANGE**: A green box with the sub-section 'Benzema régale' and a snippet of news about a match on Saturday.
- LE SAVIEZ-VOUS ?**: A blue box with a profile card for 'Karim Benzema', including his birth date (1987), age (26), club (Real Madrid), position (Attacking (football)), and height (1.87 m).
- SITES FRANCOPHONES**: A list of search results for 'Karim Benzema' with snippets and links.
- recherches associées**: A grey box listing related search terms like 'karim benzema biographie', 'karim benzema club', etc.
- VIDÉOS**: A section with video thumbnails and titles like 'Ballon d'Or - un trophée sans suspens ?' and 'FOOT - BLEUS - Benzema remis sur les rails du succès'.

At the bottom, there's a footer with 'AUTRES RECHERCHES QUI BUZZENT' and 'mises à jour toutes les heures'.

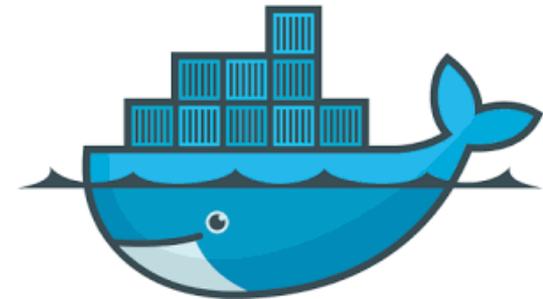
architecture d'origine : contrainte

1 bloc
200m docs



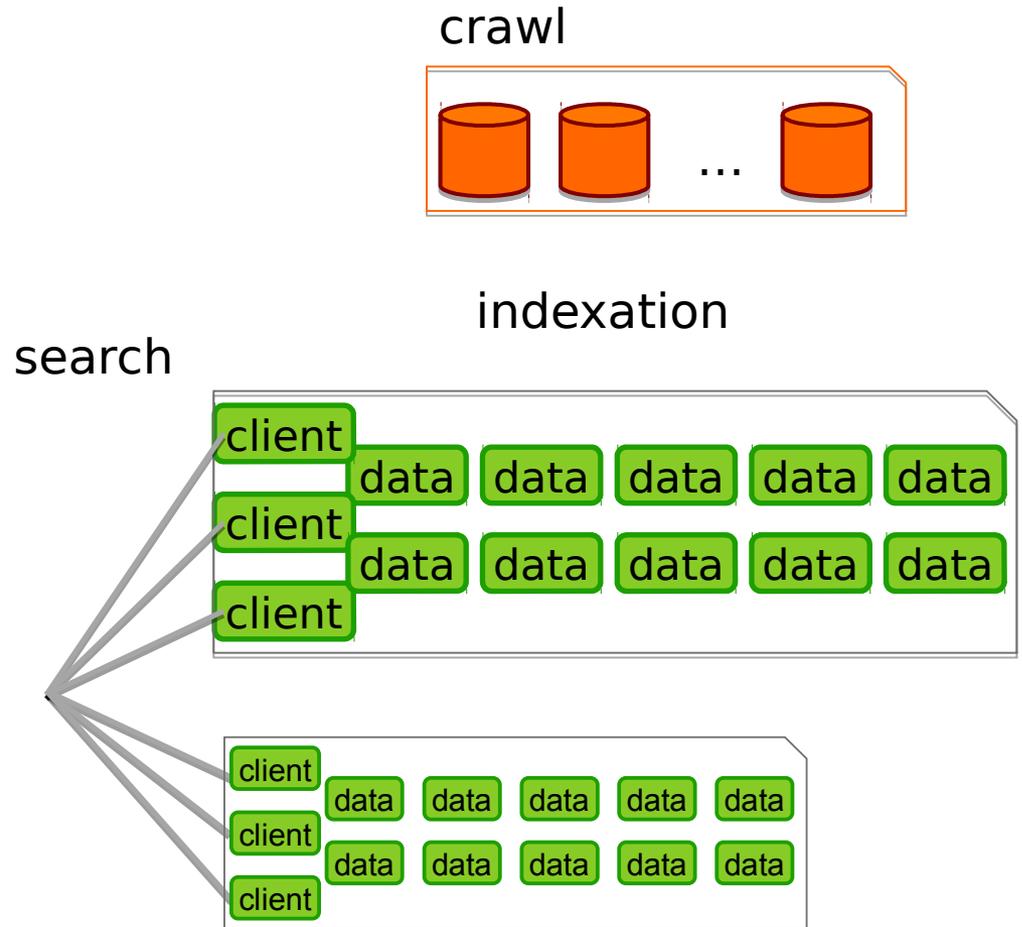
développement en production

- disposer de la volumétrie
- benchs réalistes
- isoler les applicatifs

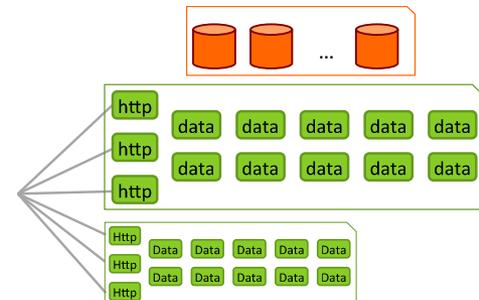
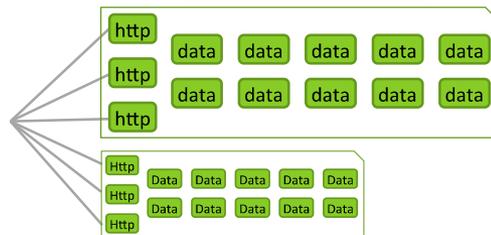
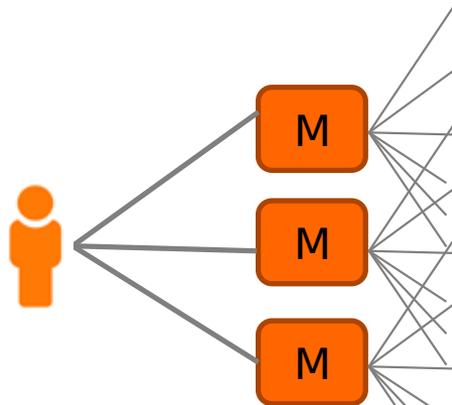


architecture de transition

1 bloc
200m docs

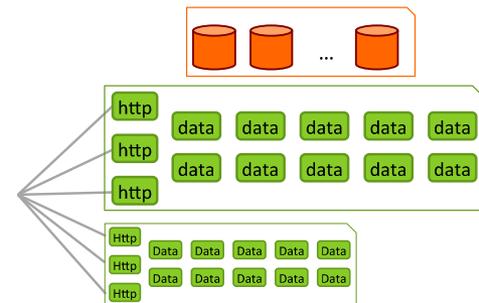
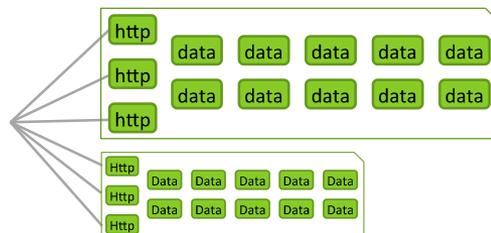
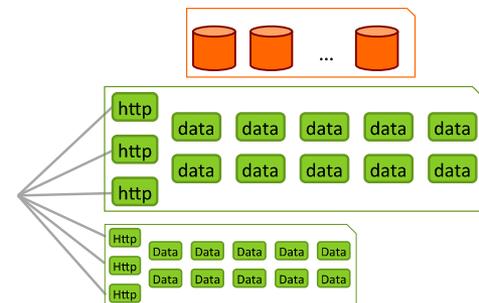
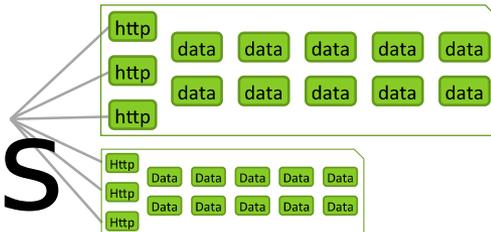


architecture de transition



6 blocs
1 200m docs

multisite



performances search

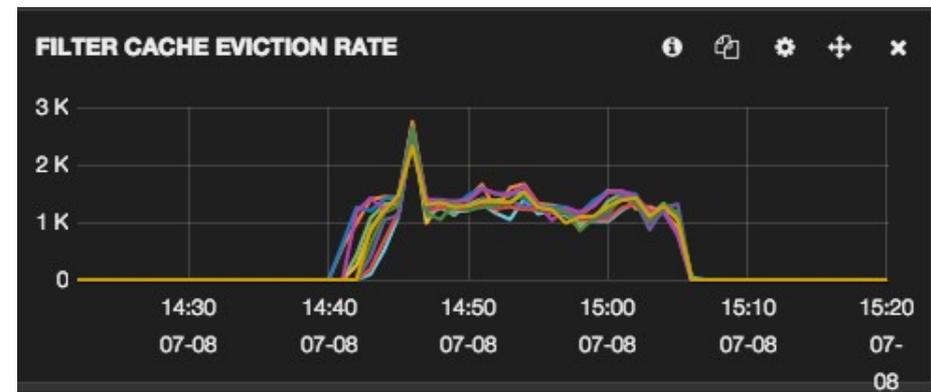
- requête JSON
 - trouver la requête le plus pertinente
 - ET la plus rapide
- filter cache size
 - 80%, 60%, 40%, 30%, 20%, 10%
 - contrôle fin
- machines
 - 10, 20
- nodes
 - 1, 2 par machine
- os.processors
 - not set, 4 8, 16
- parsing JSON
 - loader puis dispatcher
 - changement de lib C++
- shards
 - 80, 40, 20
- replica
 - 0, 1, 2
- RAM
 - 8Go, 48Go, 8Go
- agrégation
 - une requête, deux requêtes

un p'tit bout de Marvel

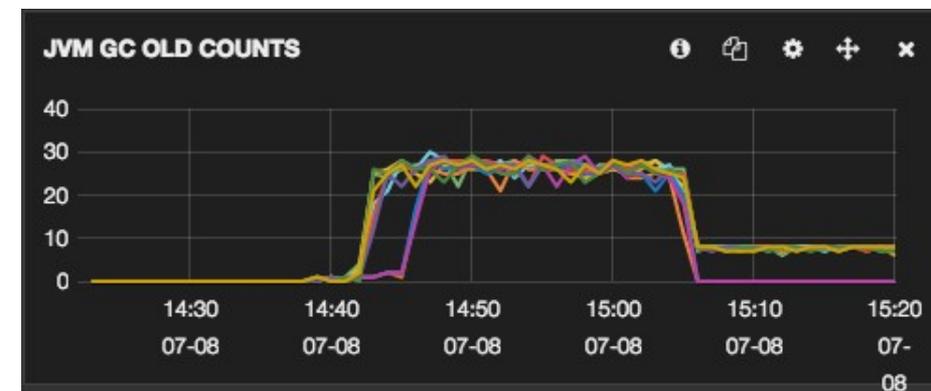
- search ops/sec



- cache plein => éviction



- éviction => GC old



bonnes pratiques

- oubliez ça !
 - testez, testez, re-restez !
- garder précieusement les résultats
 - sinon pas de comparaisons faciles

- déployez des outils de monitoring
 - plus d'infos pour mieux comprendre

- faire des petits déplacements autours des meilleurs réglages
 - pour éviter la multiplication des benches

- requête JSON
 - trouver la requête le plus pertinente
 - ET la plus rapide
- filter cache size
 - 80%, 60%, 40%, 30%, 20%, 10%
- machines
 - 10, 20
- nodes
 - 1, 2 par machine
- os.processors
 - not set, 4 8, 16
- parsing JSON
 - rapidJSON
- shards
 - 80, 40, 20
- replica
 - 0, 1, 2
- RAM
 - 8Go, 48Go, 8Go
- agrégation
 - une requête, deux requêtes

plan

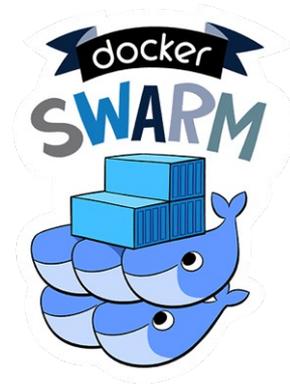
- hier
- aujourd'hui
- demain

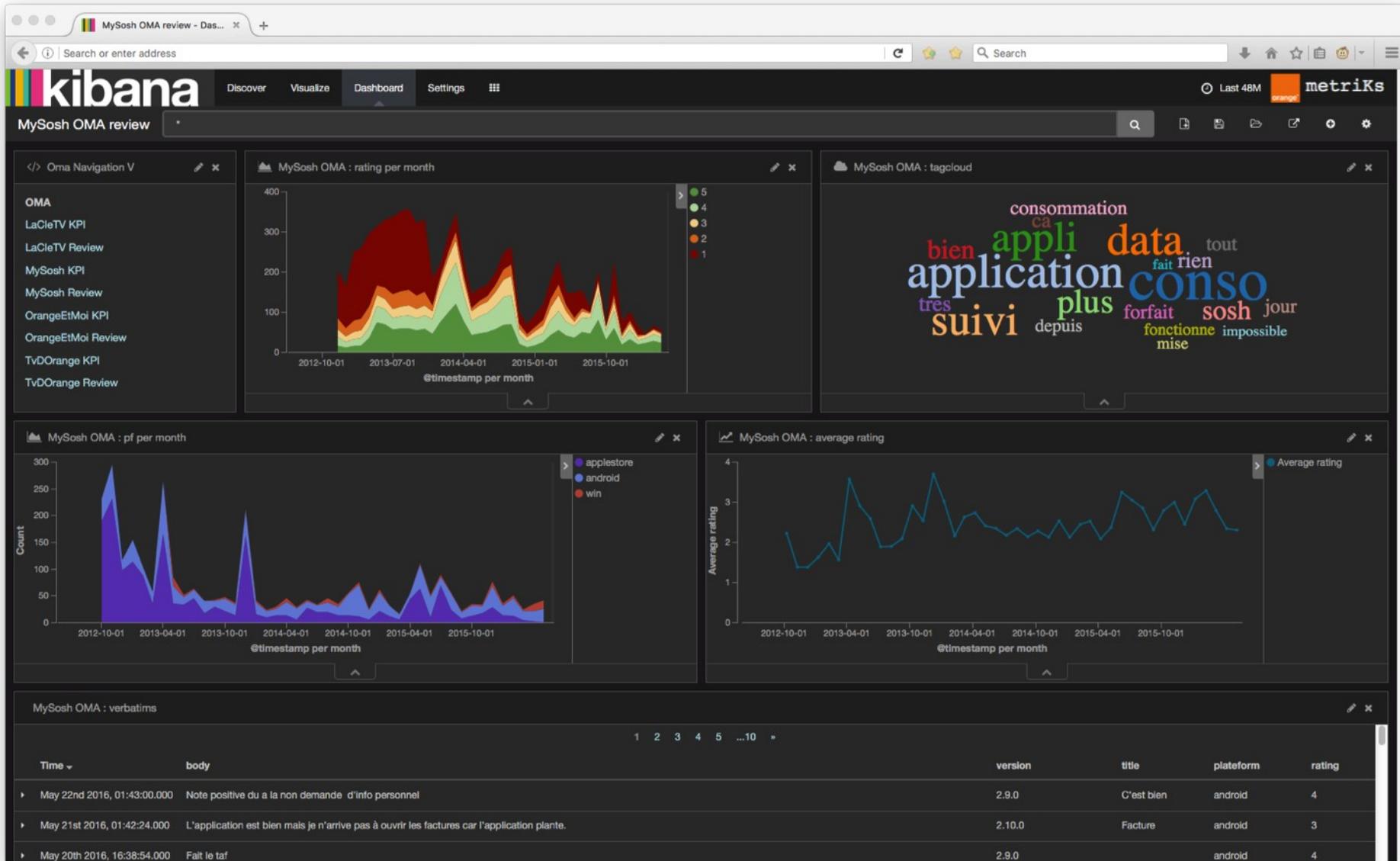
moteurs thématiques

- dockerisation
- accompagnement search
 - aide pour la création de requêtes
 - outils non-régression
- maintenance, supervision
 - suivre les versions ElasticStack

statistiques d'usage

- collecte, centralisation
 - xxxBeat
 - logstash
 - rabbitmq
- cluster élastique
 - on demand
 - séparation des données
- visualisation
 - kibana, grafana...
 - accompagner nos utilisateurs





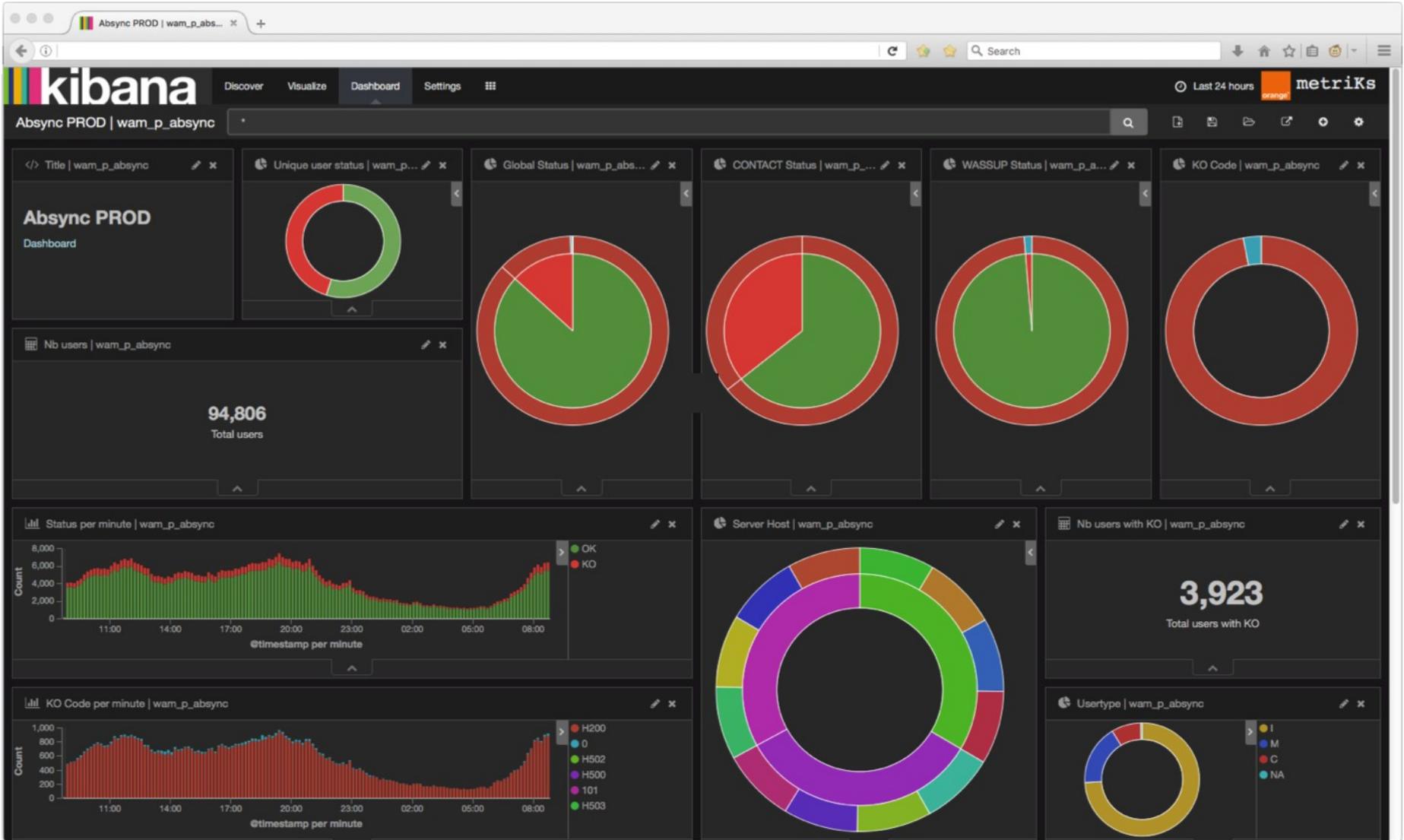
c'est l'heure du Quiz !

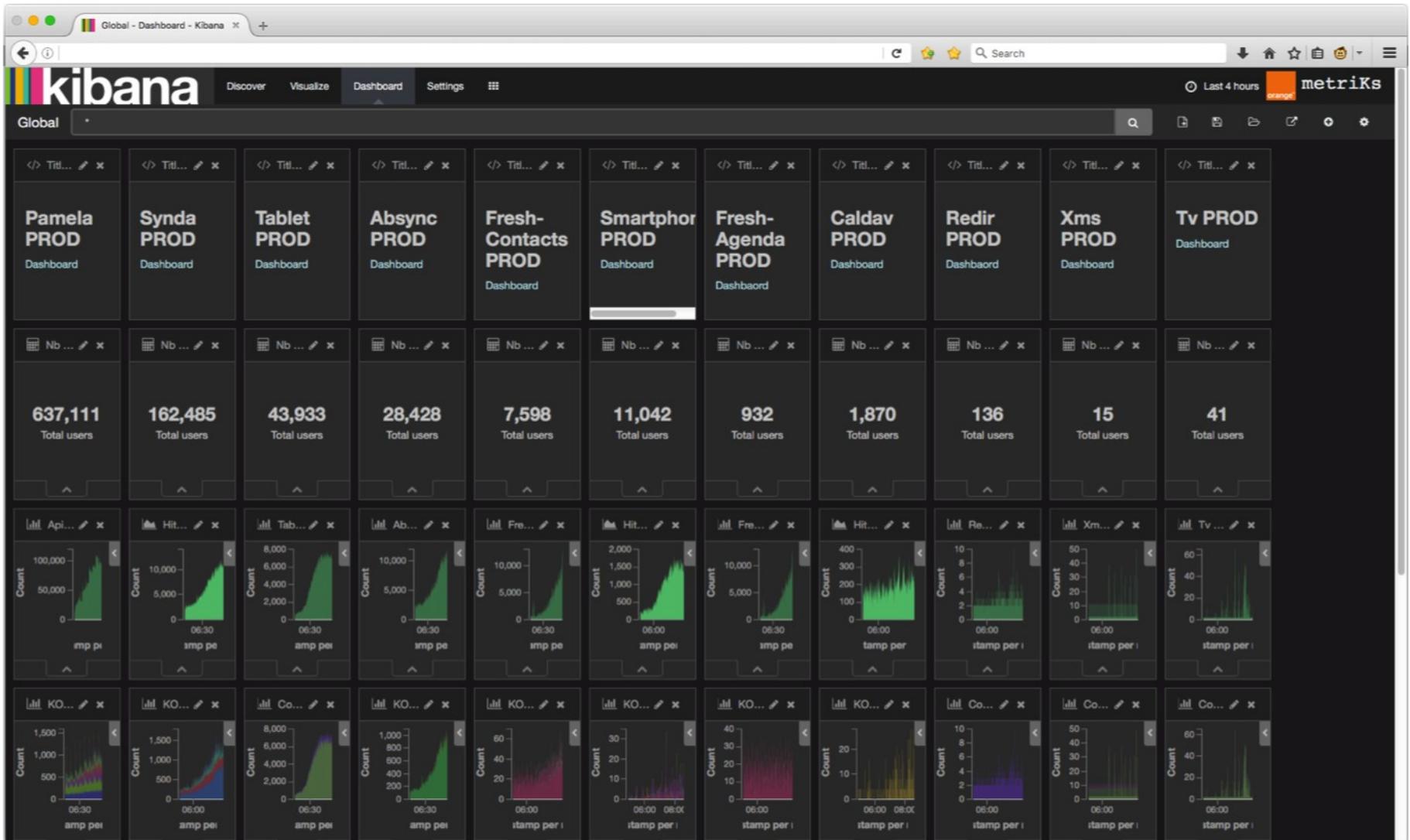
question : nos moteurs
thématiques sont en production
avec Elasticsearch depuis ?

- A) avril 2013
- B) fin 2015
- C) début 2014
- D) 1^{er} août 2010

répondez vite en tweetant sur
@TechConfQuiz







plan

- hier
- aujourd'hui
- demain

conclusion

- veille et animation technique
 - interne Orange

- calculs distribués
 - comptage unique
 - pré traitements
 - prévisions

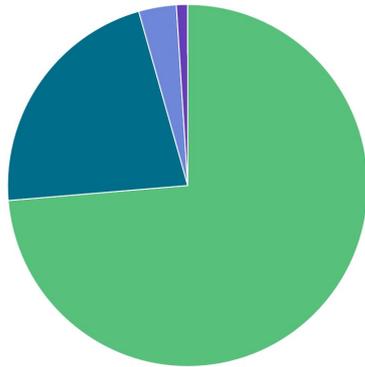
- contribuer
 - logstash, beat, kibana



conclusion

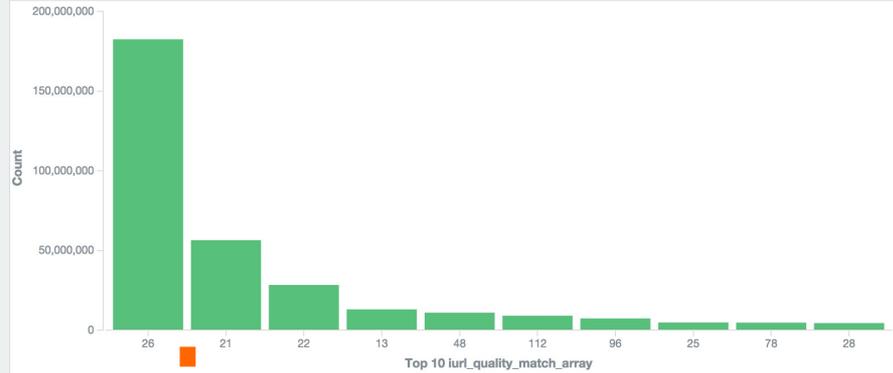
- facilité d'intégration
 - REST JSON
 - tout le temps
- flexibilité, versatilité
 - de 1000 docs à 10^{10} docs
- couverture de la stack
- documentation
- offre commerciale
- évolution rapide (trop ?)

Pourcentage bonnes/mauvaises urls



- Legend
- score_url:0
 - score_url:[1 TO 500]
 - score_url:[501 TO 1000]
 - score_url:[1001 TO 3000]
 - score_url:>3001

Top 10 iurl_quality_match_array



merci

Label

Top 10 label_match_iurl_array	Count
UNID	182,452,681
SPAM	116,694,046
USER	30,446,426
TOOL	22,334,605
SRCH	9,982,425
LANG	6,408,751
ADVS	783,554
FORM	302,648

Urls

Time	score_iurl	_id
April	5,479	http://www.lefigaro.fr/bd/2009/10/30/03014-20091030ARTFIG00023-angouleme-menaces-sur-le-festival-de-b+FIGARO+-+Culture
April	5,233	http://www.biosantebeaute.fr/2012/consommer-bio-local-et-solidaire-grace-a-biocoop/?utm_source=feedburner&utm_medium=feed
April	4,775	http://www.netwizz.net/blog/2007/05/27/473-apollo-google-analytics-widget?utm_source=feedburner&utm_medium=blog
April	4,775	http://www.sachal.fr/2007/11/28/qui-saurait-creer-un-widget-covoiturage-professionnel/?utm_source=feedburner&utm_medium=blog

Count total

1,194,486,306
Count