

# ECナビ 商品検索でのSolrの利用

株式会社 ECナビ  
システム本部 ECナビラボ  
春山 征吾

# アジェンダ

---

- ▶ ECナビでのSolrの利用
- ▶ ECナビ商品検索の概要
- ▶ システム構成
- ▶ ECNaviTokenizer
- ▶ 問題点

# ECナビでのSolrの利用

---

- ▶ 商品検索
- ▶ ソーシャルブックマークサイト Buzzurl  
での検索
- ▶ ...

# 商品検索での Solr の利用

---

以下で Solr を利用しています.

- ▶ 通常の全文検索 (EC ナビとすべての OEM)
- ▶ アイテム一覧ページ (一部の OEM を除く)
- ▶ NavicSearch API

# ECナビ商品検索の概要(続き)

---

件数 約2300万件

クエリ数/日 最大約100万

元データサイズ 約10GB

インデックスサイズ 約18GB

# 利用している Solr の機能

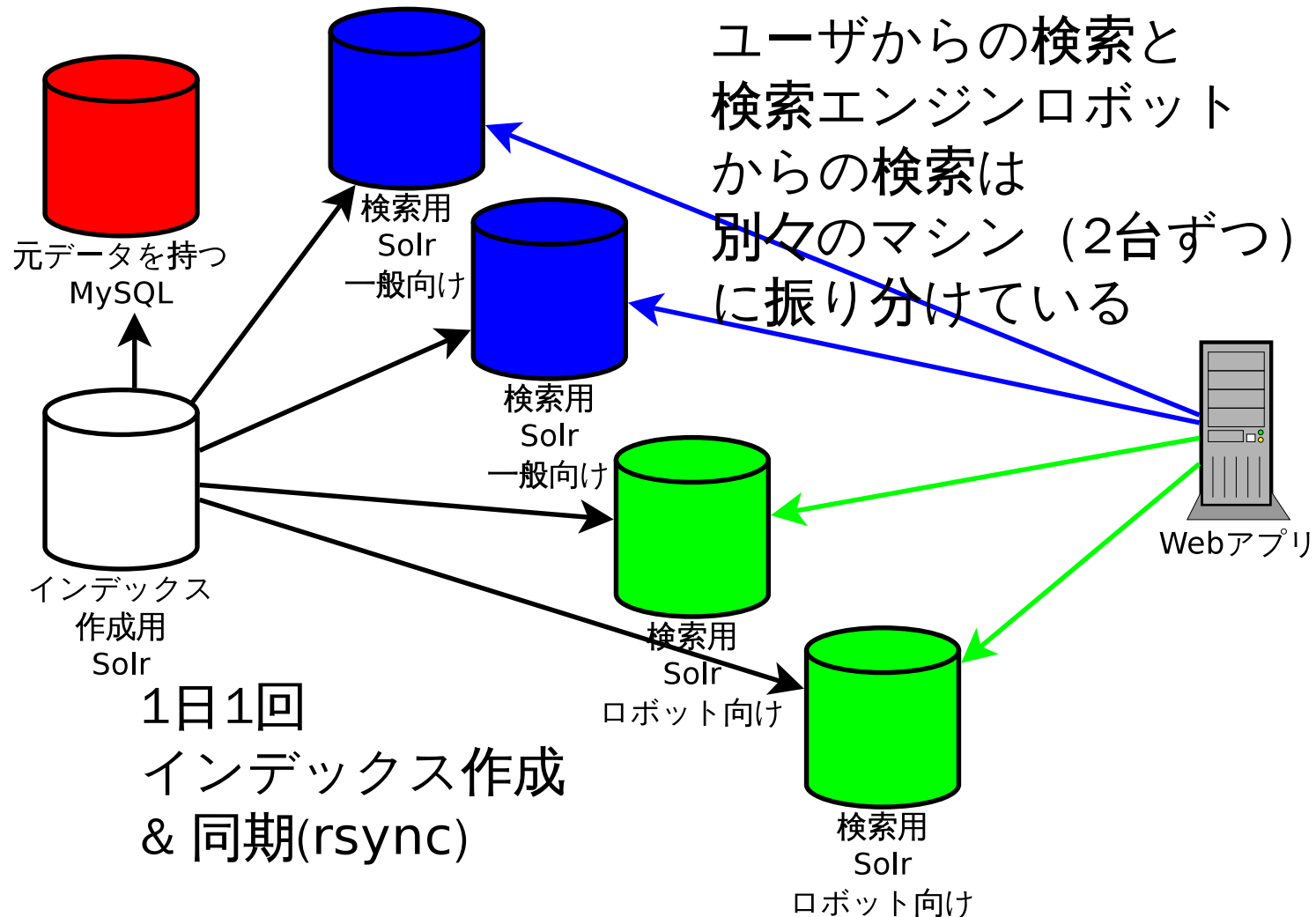
---

- ▶ DistributedSearch(shards)
- ▶ rsync でのレプリケーション
- ▶ ファセット
- ▶ 独自 Tokenizer の組み込み

# デモ

---

# システム構成





# マシンスペック

---

## インデックス作成機

メモリ 4GB

**HDD** 120GB RAID なし

## 検索機

メモリ 18GB

**HDD** 120GB RAID1+0

# ソフトウェア構成

---

- ▶ Solr 1.3
- ▶ Tomcat 6
- ▶ Sun JDK 1.6
- ▶ CentOS 4.6

# インデックス作成

---

- ▶ 1日1回インデックス作成  
rsyncでレプリケーション
- ▶ shardを3つに分けて並列に更新
- ▶ インデックス更新にかかる時間は、  
2時間

# 検索

---

- ▶ 検索エンジンのロボットからのクエリが多数
- ▶ そのため、4台のマシン中2台がロボット用2台が人間用

## 検索(続き)

---

どれくらいロボットからのクエリが多いか  
というと

ロボット 50万～80万/日

それ以外 20万弱/日

# ECNaviTokenizer

---

武田さんのCJKTokenizerを改造したものの

- ▶ 「つのだ ひろ」「漫 画太郎」といった語句の検索のために「 」などの文字があってもなくても同じTokenが出力されるようにした

# ECナビ独自のTokenizer(続き)

---

- ▶ 商品の型番によくある「cyber-shot」「cybershot」などといったハイフンの表記揺れに対応するため、「cyber-shot」から「cyber」「shot」「cybershot」というTokenを生成するようにした。(検索時には、すべてのハイフンは取り除いてクエリが送られる)
- ▶ いわゆる半角カナを全角カナに正規化するようにした

# 問題点

---

- ▶ インデックスがメモリに乗らない大きさ
- ▶ 検索エンジンロボット由来のクエリにはキャッシュの効果がない
- ▶ Disk Accessが発生      HDDがボトルネック



というわけでLTに続く

---