

# Pipelining your music



**<whoami>**

**Jyrki Pulliainen**

Content team



# Spotifier since January

Pythonista since 2.3



@nailor



`</whoami>`



# Spotify?



Spotify Premium - Linux Preview

File Edit View Playback Develop Help

numminen boogie

MAIN

- What's New
- People
- Inbox
- Play Queue
- Devices

APPS

- App Finder
- Top Lists
- Radio
- Classify
- The Guardian
- Last.fm
- Moodagent
- Pitchfork
- Songkick Concerts
- Soundrop

COLLECTION

- Library
- Local Files
- Downloads

CHRIS BROWN

LYSSNA NU

Innehåller hitsen Turn Up The Music och Don't Wake Me Up  
Gästas av B.O.B., Snoop Dogg, Wiz Khalifa och Nas

THE VEON TTES B SER A TON

DJ-Kicks  
Digitalism

DIGITALISM  
RECOMMENDED  
DJ-KICKS

ADRIAN LUX  
FEATURING DANTE

TURNING

Featured Spotify Apps

1. TuneWiki  
Get synced lyrics and sing along to your favorite songs!

2. Tunigo  
Ett enkelt och underhållande sätt att upptäcka ny musik.

3. Digster  
Playlists from your favorite artists and for all occasions!

4. Filtr  
Playlist generator with a social twist.

For Silence  
Flying Saucer Attack

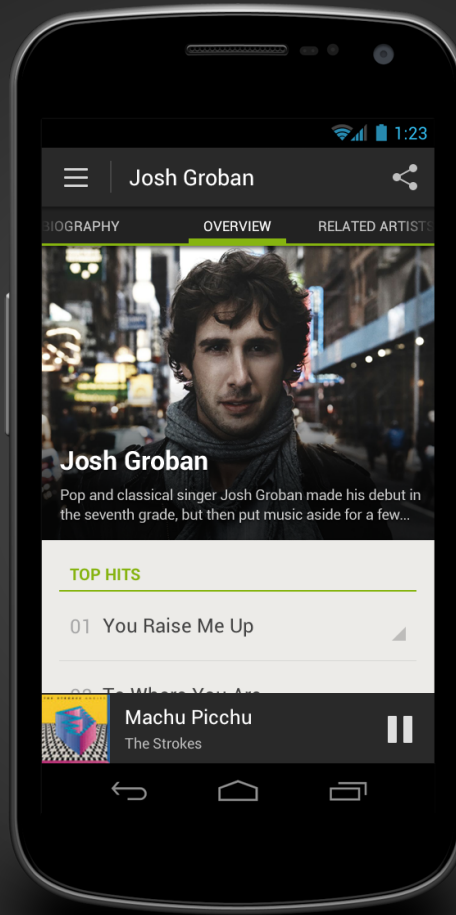
Jyrki Pulliainen

- Aino Kopu
- Anders Åstrand  
All I Ever Wanted by Adria...
- Jarno Tikka  
The Gift Of Guilt by Gojra
- Johanna Henriksso
- Johannes Dahlström
- Kalle Rindell
- Lasse Saarinen
- mharju  
Holding Out For A Hero by
- Mikko Harju
- Richard Bengtsson  
listened to Villi You Follo...
- Anders Åstrand listened to All I Ever Wanted by...
- Pierre Carrier listened to Never Mind by Infected...
- Ville Kaisla listened to Perculator Meme (Lazer...
- Mats Linander listened to Skinny Love by Bon Iver
- Anders Åstrand listened to Clean Gloves, Dirty...
- Mats Linander listened to Lump Sum by Bon Iver
- Richard Bengtsson listened to Hollow Crown...
- Ville Kaisla listened to Busy Crack by Teeth
- Anders Åstrand listened to Too Close by Alex Clare
- Mats Linander listened to Flame by Bon Iver
- Richard Bengtsson

0:09 7:39







# Number fun

- 10M monthly active users
- 18M tracks
- 100 years of music
- 20k added **every** day



# The Music Pipeline

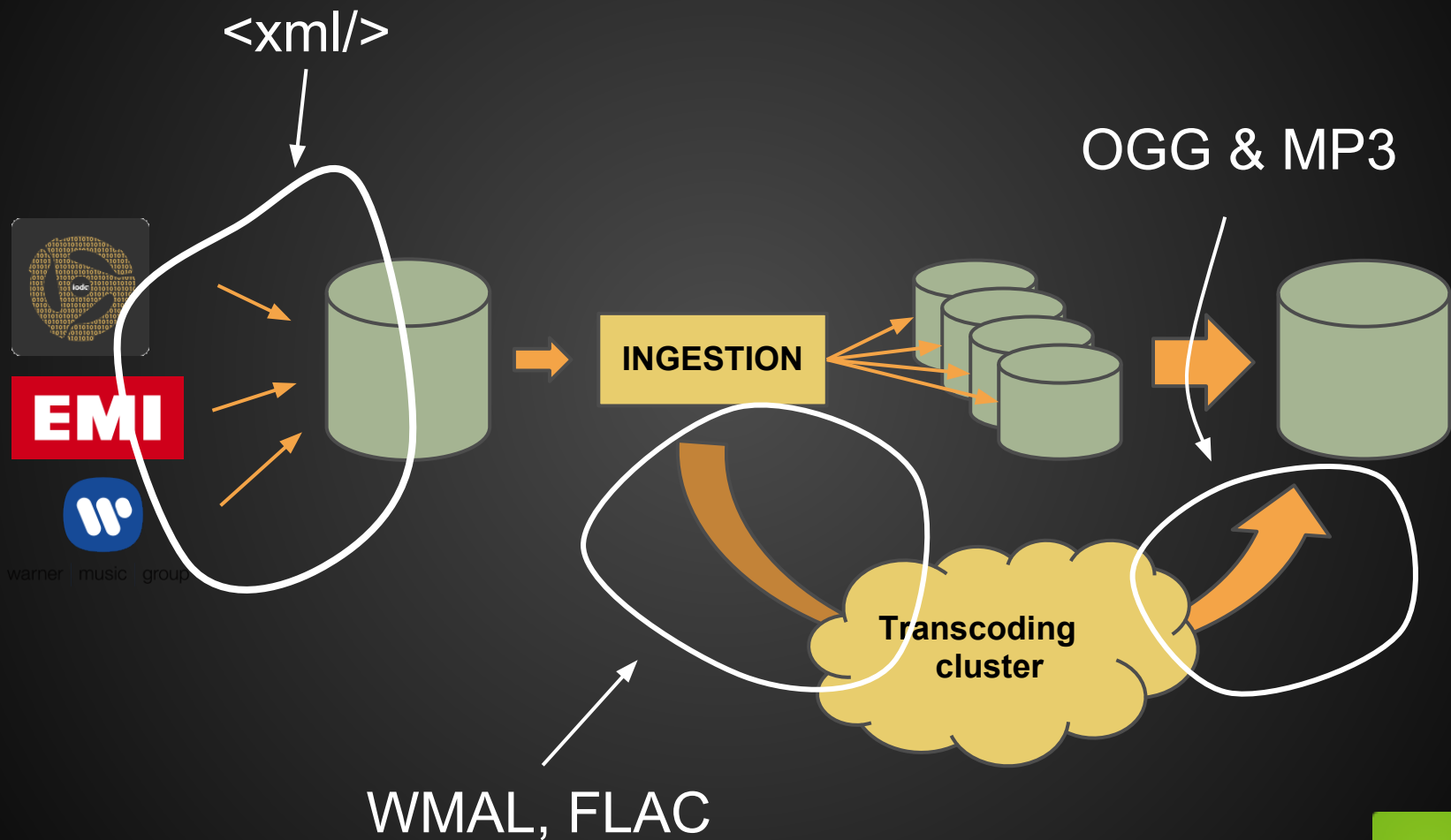


# ANTI HIPSTER



# STACK





**100s of TBs of data**

**Load of deliveries daily**

**Malformed data every day**



A close-up photograph of a hand holding a hammer. The hammer has a wooden handle and a metal head. A semi-transparent grey rectangular box with a white border is overlaid diagonally across the image, containing the text "Right tools for the right job" in white, bold, sans-serif font.

**Right tools for the right job**

# XPath extensions



```
>>> def formerlify(_, name):  
...     return 'The artist formerly known as %s' %name  
>>> #Namespace stuff  
>>> from lxml import etree  
>>> ns = etree.FunctionNamespace('http://my.org/myfunctions')  
>>> ns['hello'] = hello  
>>> ns.prefix = 'f'  
>>> root = etree.XML('<a><b>Prince</b></a>')  
>>> print(root.xpath('f:hello(string(b))'))  
... The artist formerly known as Prince
```



# Fun(?) facts

- 10 different XML formats
  - Majors vs our own (indies)
  - One industry "standard"

Biggest XML 3.3M lines (350MB)

- Bible apparently fits in 3MB of XML
- lxml ftw



```
>>> min(timeit.repeat('etree.parse("huge.xml")', setup="from
lxml import etree", number=1, repeat=5))
2.309144973754883
```

```
>>> min(timeit.repeat('etree.parse("huge.xml")', setup="from
xml.etree import cElementTree as etree", number=1, repeat=5))
3.0681779384613037
```

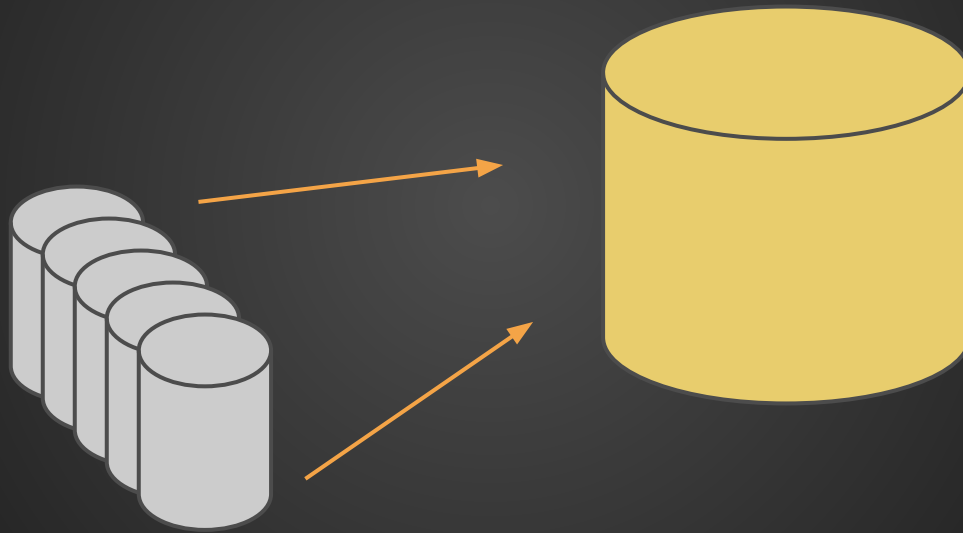
```
>>> min(timeit.timeit('etree.parse("huge.xml")', setup="from
xml.etree import ElementTree as etree", repeat=5, number=1))
Killed
```

```
>>> # (with PyPy 1.9)
```

```
>>> min(timeit.repeat('etree.parse("huge.xml")', setup="from
xml.etree import ElementTree as etree, number=1, repeat=5))
23.186518907546997
```



# Merging



# Fun(?) facts

- Artists don't have any global or even label specific IDs
  - Multiple artists with same name
  - Even spelling differs inside a single label
- Multiple versions of the same album
- Enormous search space!
  - $(18 * 10^{**6}) ** 2 ==$  huge number



Machine learning!

Insufficient  
data

**Tip: Reduce search space!**

```
>>> from unicodedata import normalize
>>> key = ''.join(normalize('NFD', char)[0].lower() for char
in title)[5]
```

**Side note: Levenshtein is expensive  
=> use other edit distances too**

(or use PyPy, 4x speed increase ftw)







# Pro-tips

- If data is relational, use relational database (duh)
- Don't over-normalize yourself, BCNF is rarely beneficial
- Weight between denormalize vs. moar indices
- Let the DB do the hard lifting, query planner is your friend!





- Asynchronous!
- RabbitMQ + amqp-lib
  - One master, 49 slaves
- Isilon storage => 8Gbit/s throughput!

1M / DAY

# Index building

....with Java



# Why not Python?

- Not powerful enough for computationally intensive stuff
- We use Lucene for Search, so Java is a natural choice

...but I'd like to try PyPy here.



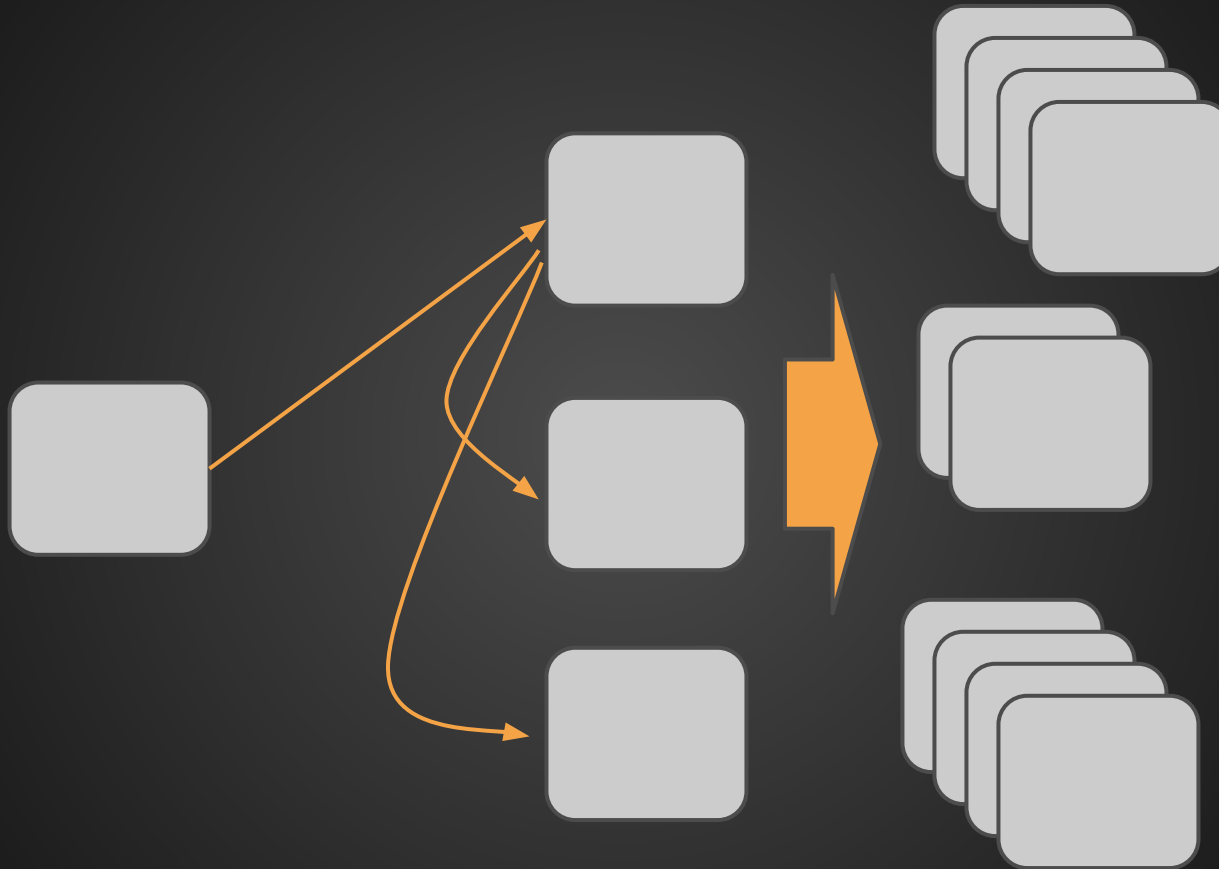
# The Music Distribution



**TTL one day**  
time to live (laiv)



# Publishing an index





**SCPing around, moving  
hundreds of GBs daily**

**Future == BitTorrent**

...not totally free of issues either



# Index format?

Read only K/V (mostly)



**Keep your eye on it**



**Mind the speed!**



# Experiment



**Ditch your code**



# Thank you!

[spoti.fi/ep\\_2012](https://open.spotify.com/playlist/spoti.fi/ep_2012)

