

Unified Pipeline Architecture: The Evolution of Data Processing at Spotify.

Erin Palmer
Applied Data Scientist @ Spotify



What's Spotify, and What Data Do We Process in Creator?



The Beatles

8,961 listeners now

[Your Artists](#)

[Help](#)

[Sign Out](#)

Monthly Listeners

Daily Listeners Fans



7.2M Monthly Listeners ?

Down 8.8k yesterday



Compare to

🔍 2 more artists



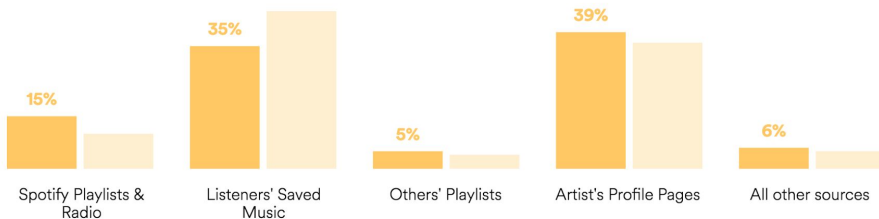
AUDIENCE

SONGS

PLAYLISTS

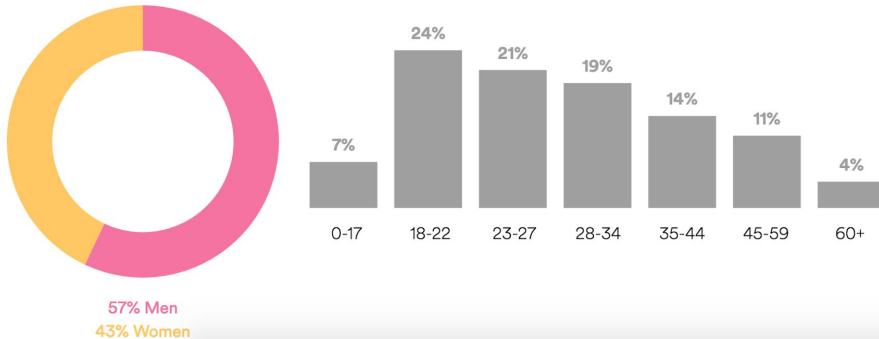
How They Listen on Spotify

Monthly Listeners Fans



Who They Are

Monthly Listeners Fans



What's Our Primary Source?

EndSong

- Log record containing:
 - trackGid
 - userId
 - Ms Played
 - Time started
 - Play Context
 - Location
 - Etc
- From it we can derive most of the useful information for listening history

Secondary Sources

- User Meta Data
- Playlist Snapshot
- Artist/Track Meta Data
- etc

Current Architecture

Characteristics

- Independent code paths compute necessary endpoints
- Each endpoint has 3-6 intermediate computation stages
- Intermediate Data is not reusable

Current Architecture

Advantages

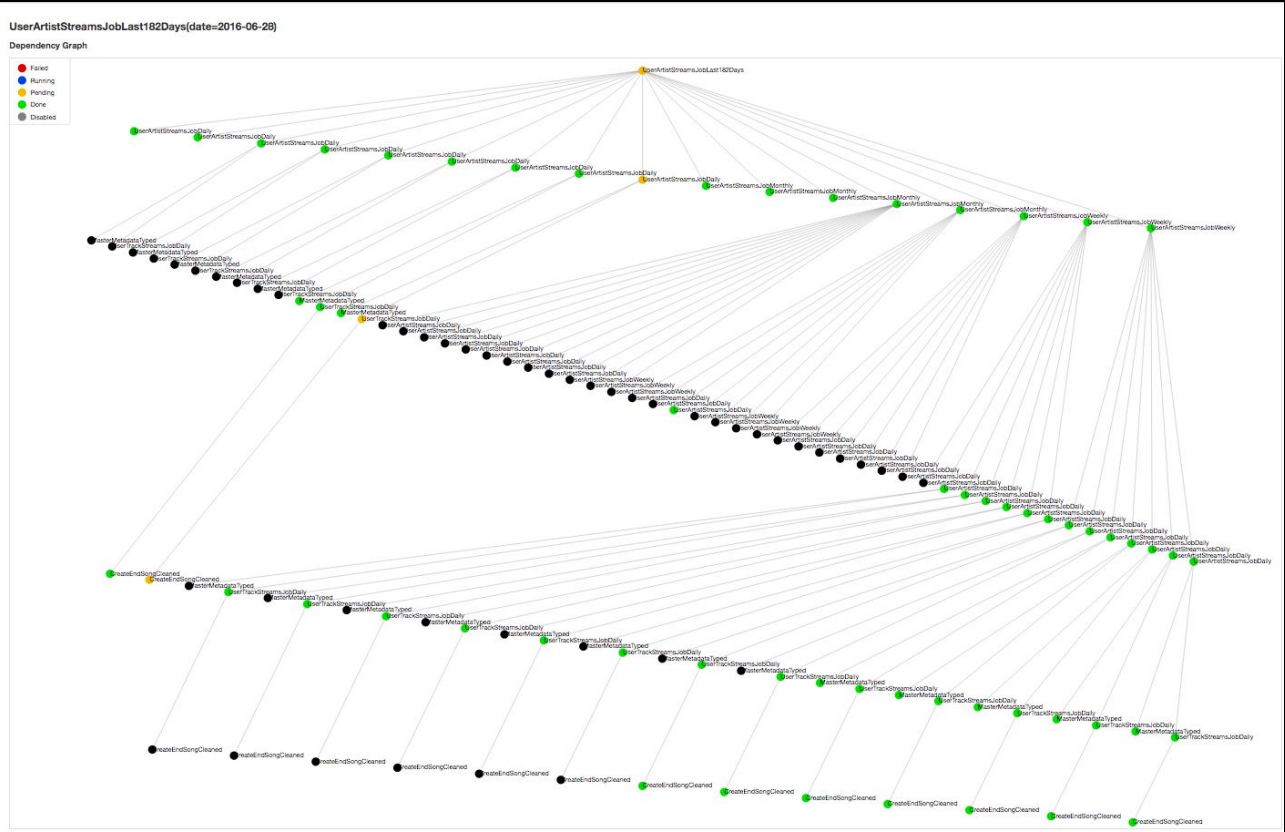
- Very Flexible
- New pipelines can be easily built
- Adaptable to constantly changing requirements

Current Architecture

Issues

- Redundant Dependencies
- Multi-layer dependencies cause cascading effects if there are any issues/delays
- Inconsistencies

Redundant Dependencies





Redundant Dependencies

Source	Target
ArtistRankHistoryDailyJob	ArtistRankHistoryDailyJob
ArtistRecordingMetricsJob	ArtistRecordingMetricsJob
ArtistStreamsAggregation	ArtistStreamsAggregation
ArtistStreamsDaily	ArtistStreamsDaily
BingedJob	BingedJob
BingerSegmentSummaryAggregation	BingerSegmentSummaryAggregation
CityArtistStreamsAggregationJob	CityArtistStreamsAggregationJob
CityListenerArtistAggregation	CityListenerArtistAggregation
CityListenerRecordingAggregation	CityListenerRecordingAggregation
CityStreamsJobDaily	CityStreamsJobDaily
CityToplistArtistFanAggregation	CityToplistArtistFanAggregation
CityToplistArtistListenerAggregation	CityToplistArtistListenerAggregation
CityToplistRecordingAggregation	CityToplistRecordingAggregation
CityTrackStreamsAggregationJob	CityTrackStreamsAggregationJob
coredata.CreateEndSongCleaned	coredata.CreateEndSongCleaned
creator.ArtistStreamExportJobESAlias	creator.ArtistStreamExportJobESAlias
creator.BingerSegmentSummaryExportJobESAlias	creator.BingerSegmentSummaryExportJobESAlias
creator.CityFanExportJobESAlias	creator.CityFanExportJobESAlias
creator.CityListenerExportJobESAlias	creator.CityListenerExportJobESAlias
creator.FanCountryExportJobESAlias	creator.FanCountryExportJobESAlias
creator.FanDemographicsExportJobESAlias	creator.FanDemographicsExportJobESAlias
creator.FanStreamSourceExportJobESAlias	creator.FanStreamSourceExportJobESAlias
creator.FanTimelineExportJobESAlias	creator.FanTimelineExportJobESAlias
creator.GeoLocationSuggestionsESAliasSwapVerify	creator.GeoLocationSuggestionsESAliasSwapVerify
creator.ListenerDailyTimelineExportJobESAlias	creator.ListenerDailyTimelineExportJobESAlias
creator.Listener28DayCountryExportJobESAlias	creator.Listener28DayCountryExportJobESAlias
creator.Listener28DayDemographicsExportJobESAlias	creator.Listener28DayDemographicsExportJobESAlias
creator.Listener28DayStreamSourceExportJobESAlias	creator.Listener28DayStreamSourceExportJobESAlias
creator.Listener28DayTimelineExportJobESAlias	creator.Listener28DayTimelineExportJobESAlias
creator.ListenerFanVennDiagramExportJobESAlias	creator.ListenerFanVennDiagramExportJobESAlias
creator.LoyaltySegmentSummaryExportJobESAlias	creator.LoyaltySegmentSummaryExportJobESAlias
creator.RegularSegmentSummaryExportJobESAlias	creator.RegularSegmentSummaryExportJobESAlias
creator.SimilarArtistsExportJobESAlias	creator.SimilarArtistsExportJobESAlias
creator.SongCityMetricsExportJob	creator.SongCityMetricsExportJob
creator.SongListeningMetricsExportJob	creator.SongListeningMetricsExportJob
creator.SongCountryExportJob	creator.SongCountryExportJob
creator.SongSourceExportJob	creator.SongSourceExportJob
creator.SongTimelineExportJob	creator.SongTimelineExportJob
creator.TopPlaylistsDayExportJobESAlias	creator.TopPlaylistsDayExportJobESAlias
creator.TopPlaylists28DayExportJobESAlias	creator.TopPlaylists28DayExportJobESAlias
creator.TopPlaylists7DayExportJobESAlias	creator.TopPlaylists7DayExportJobESAlias
creator.TopSongPlaylistsDayExportJob	creator.TopSongPlaylistsDayExportJob
creator.TopSongPlaylists28DayExportJob	creator.TopSongPlaylists28DayExportJob
creator.TopSongPlaylists7DayExportJob	creator.TopSongPlaylists7DayExportJob
creator.TopSongExportJob	creator.TopSongExportJob
FanSegmentAgeGenderAggregation	FanSegmentAgeGenderAggregation
FanSegmentCountryAggregation	FanSegmentCountryAggregation
FanUnionSegmentDemographicsByCountryJob	FanUnionSegmentDemographicsByCountryJob
FanUnionSegmentSizeAggregation	FanUnionSegmentSizeAggregation
FanUnionStreamSourcesAggregation	FanUnionStreamSourcesAggregation
ListenerFanIntersectionSizeAggregation	ListenerFanIntersectionSizeAggregation
ListenerJobDaily	ListenerJobDaily
ListenerJobLast28Days	ListenerJobLast28Days
ListenerSegmentAgeGenderAggregation	ListenerSegmentAgeGenderAggregation
ListenerSegmentCountryAggregation	ListenerSegmentCountryAggregation
ListenerSegmentDemographicsByCountryJob	ListenerSegmentDemographicsByCountryJob
ListenerSegmentSizeAggregationDaily	ListenerSegmentSizeAggregationDaily
ListenerSegmentSizeAggregationLast28Days	ListenerSegmentSizeAggregationLast28Days
ListenerSegmentSizeAggregationLast7Days	ListenerSegmentSizeAggregationLast7Days
ListenerStreamSourcesAggregation	ListenerStreamSourcesAggregation
LoyaltyJob	LoyaltyJob
LoyaltySegmentSummaryAggregation	LoyaltySegmentSummaryAggregation
RecordingCountryAggregation	RecordingCountryAggregation
RecordingStreamsListenersJob	RecordingStreamsListenersJob
RegularFanJob	RegularFanJob
RegularSegmentSummaryAggregation	RegularSegmentSummaryAggregation
SaverFanJob	SaverFanJob
SegmentRecordingStreamSourcesAggregationJob	SegmentRecordingStreamSourcesAggregationJob
SegmentTrackStreamSourcesAggregationJob	SegmentTrackStreamSourcesAggregationJob
SimilarArtistsAggregation	SimilarArtistsAggregation
SongCityMetricsRecordingAggregationJob	SongCityMetricsRecordingAggregationJob
TopPlaylistsAggregationDaily	TopPlaylistsAggregationDaily
TopPlaylistsAggregationLast28Days	TopPlaylistsAggregationLast28Days
TopPlaylistsAggregationLast7Days	TopPlaylistsAggregationLast7Days
TopPlaylistsForRecordingAggregationDaily	TopPlaylistsForRecordingAggregationDaily
TopPlaylistsForRecordingAggregationLast28Days	TopPlaylistsForRecordingAggregationLast28Days
TopPlaylistsForRecordingAggregationLast7Days	TopPlaylistsForRecordingAggregationLast7Days
TopRecordingAggregation	TopRecordingAggregation
TrackToArtist	TrackToArtist
TrackTotalSavesJobDaily	TrackTotalSavesJobDaily
TrackUserCollectionSavesJobDaily	TrackUserCollectionSavesJobDaily
TrackUserPlaylistSavesJobDaily	TrackUserPlaylistSavesJobDaily
TrackUserTotalSavesJobDaily	TrackUserTotalSavesJobDaily
UserArtistCollectionSavesJobDaily	UserArtistCollectionSavesJobDaily
UserArtistFollowsJobDaily	UserArtistFollowsJobDaily
UserArtistPlaylistSavesJobDaily	UserArtistPlaylistSavesJobDaily
UserArtistStreamsJobDaily	UserArtistStreamsJobDaily
UserArtistStreamsJobLast82Days	UserArtistStreamsJobLast82Days
UserArtistStreamsJobLast28Days	UserArtistStreamsJobLast28Days
UserArtistStreamsJobLast7Days	UserArtistStreamsJobLast7Days
UserArtistStreamsJobMonthly	UserArtistStreamsJobMonthly
UserArtistStreamsJobWeekly	UserArtistStreamsJobWeekly
UserPlaylistTrackJob	UserPlaylistTrackJob
UserTrackArtistUserCaseJobDaily	UserTrackArtistUserCaseJobDaily
UserTrackMetricsJobDaily	UserTrackMetricsJobDaily
UserTrackRecordingArtistJoinedJobLast28Days	UserTrackRecordingArtistJoinedJobLast28Days
UserTrackSavesJobDaily	UserTrackSavesJobDaily
UserTrackStreamsJobDaily	UserTrackStreamsJobDaily



Data Inconsistencies

Alexander Forselius

Live listener count

12.2k
Monthly Listeners
Up 1.7k

12
Fans
Up 10

AUDIENCE SONGS **PLAYLISTS**

Top 1 [?] LAST 28 DAYS ▾

LISTENERS STREAMS

SPA Treatment	9.5k	15.2k ↗
---------------	------	---------

Alexander Forselius | Swimm...

3k
Daily Listeners
Up 248

3.4k
Daily Streams
Up 270

AUDIENCE **PLAYLISTS**

Top 1 [?] LAST 28 DAYS ▾

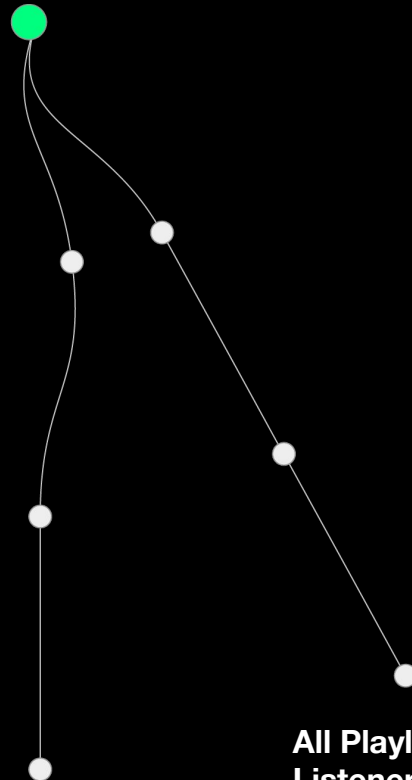
LISTENERS STREAMS

SPA Treatment	8.1k	37.2k ↗
---------------	------	---------



Current Architecture

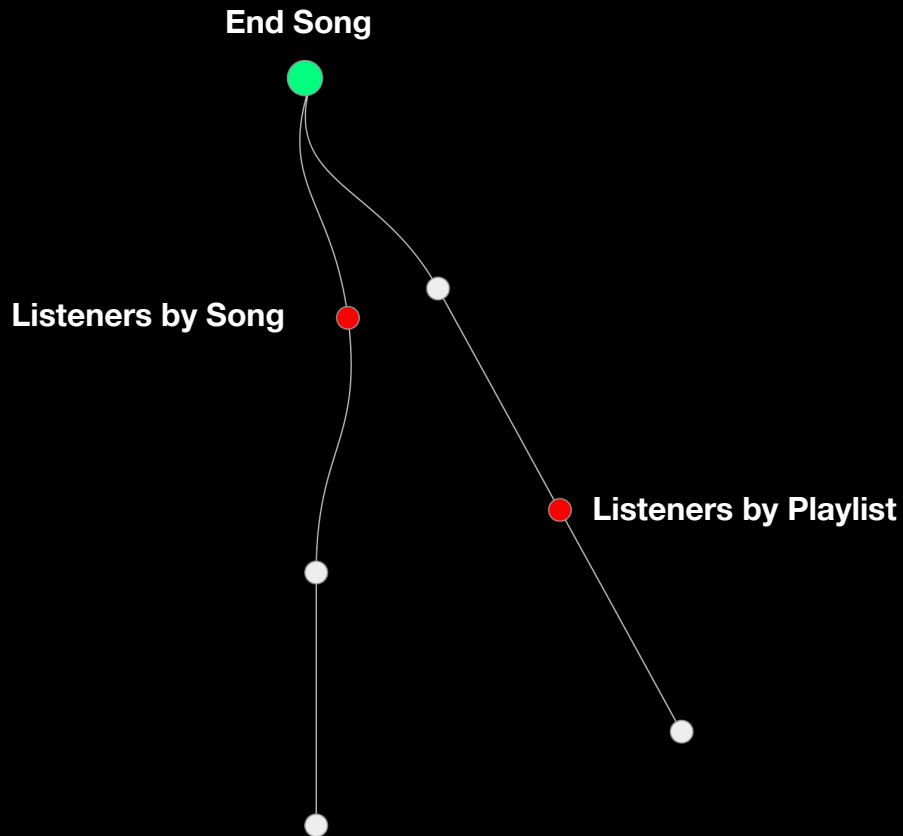
End Song



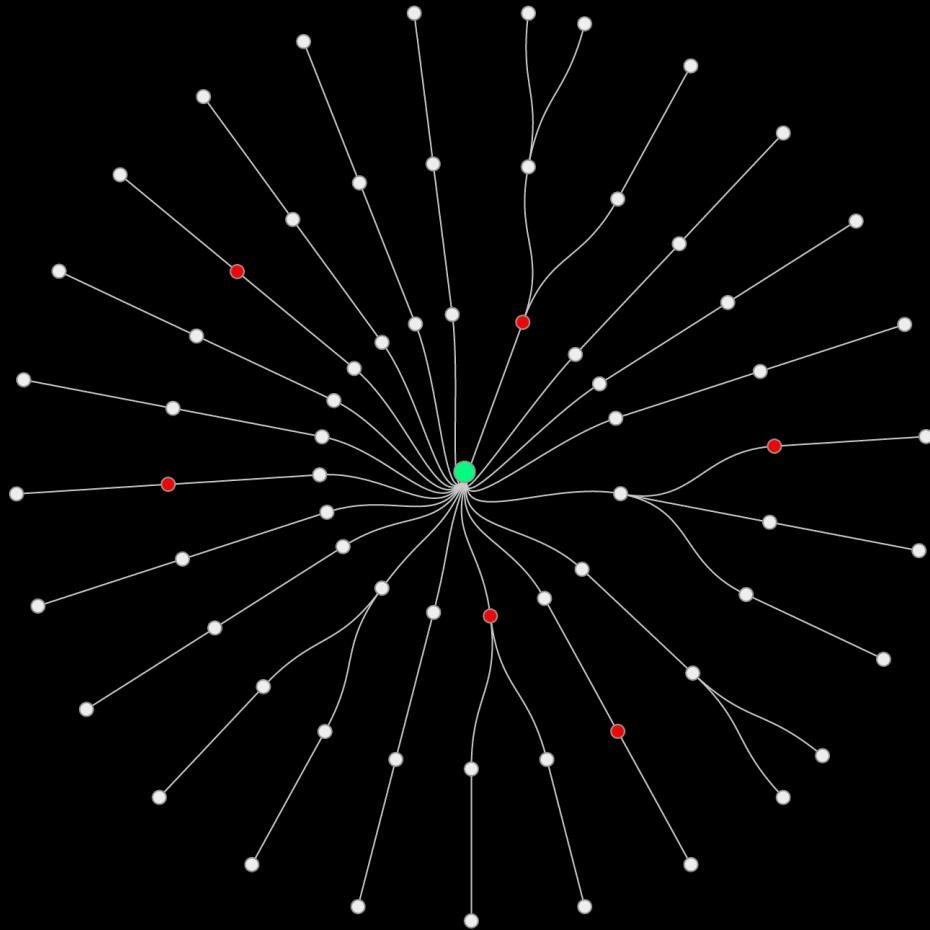
All Song
Listeners

All Playlist Song
Listeners

Duplicative Data Sets



Duplicative Data Sets

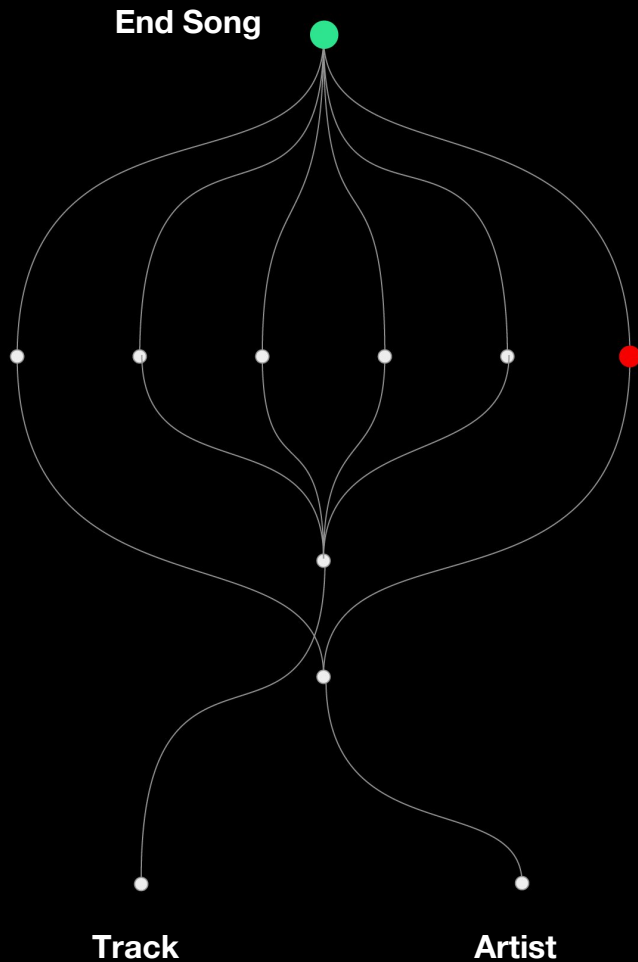


Considerations

1. Can we only read each source once?
2. How can we reduce the computation time?
3. Can we join in all the secondary sources as they comes, rather than waiting until the end of the day?
4. Retain flexibility to add new datasets or new fields



New Architecture



Architecture Components

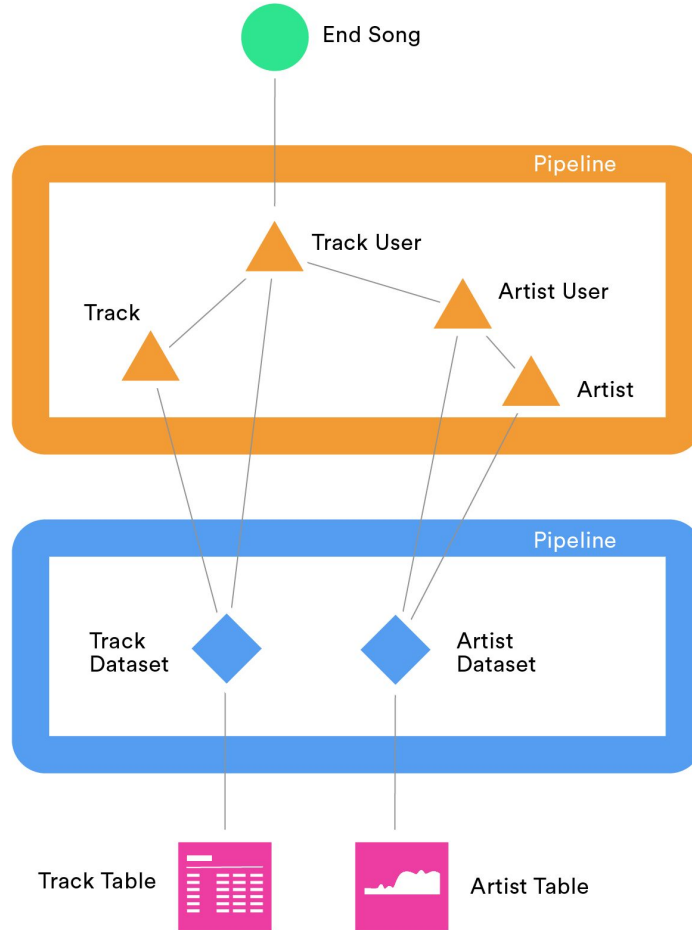
- Sources (EndSong)
- Entities
- Datasets
- Exports (Entries in a DB)

Raw Data

Entities

Datasets

Exports



Spotify
21

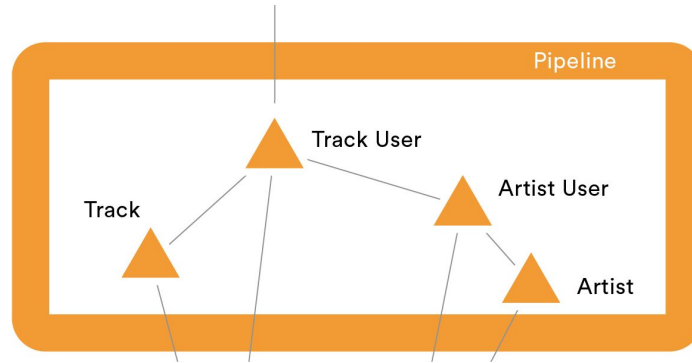
Raw Data



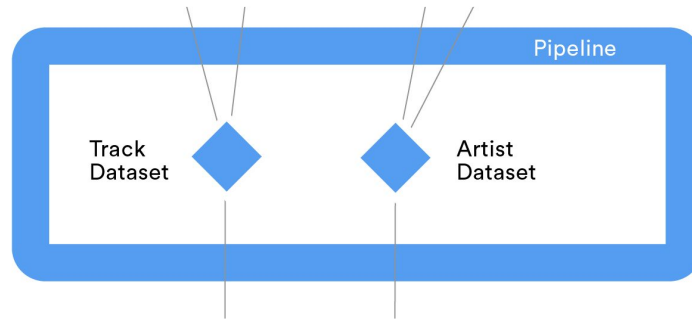
End Song



Entities



Datasets



Exports

Track Table



Artist Table

Pipeline Architecture

- Google Computing Engine
 - Google Cloud Store
 - Dataflow
 - Datastore
- Framework: Scala / Scio
 - Built on top of Google's Dataflow
- Schema: Protobufs
 - Easy Iteration
 - Built in versioning

Scio

Ecclesiastical Latin IPA: /'ʃi.o/, ['ʃi:.o], ['ʃi.i̯o]

Verb: I can, know, understand, have knowledge.

Google Dataflow with Scala => Scio

Dataflow

- Hosted, fully managed, no ops
- GCP ecosystem - BigQuery, Bigtable, Datastore, Pubsub
- Unified batch and streaming model

Scala

- High Level DSL
- Functional Programming is a natural fit for data
- Numerical Libraries: Breeze, Algebird



Example: Word Count

```
val sc = ScioContext()
sc.textFile("shakespeare.txt")
  .flatMap { _
    .split("[^a-zA-Z']+")
    .filter(_.nonEmpty)
  }
  .countByValue
  .saveAsTextFile("wordcount.txt")
sc.close()
```

Protocol Buffers

Why Protobufs?

- Land themselves to high level schema organization
- Easy to read and manipulate
- Features
 - Allows repeated fields
 - Field numerical tags allow for schema compatibility
 - Can be built to be arbitrarily large
- Protobufs are compact when stored: they are serialized into binary format

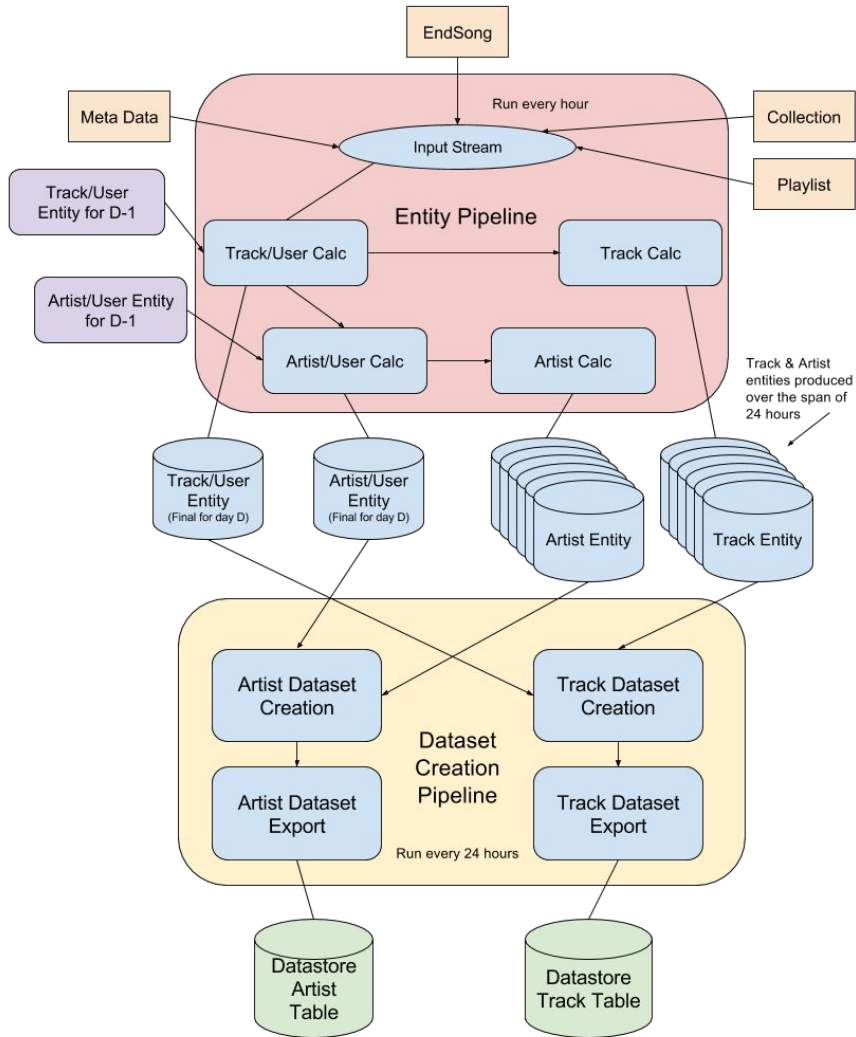
Example

```
// Next ID: 6
message AggregateKeyPB {
  enum Type {
    TRACK = 0;
    ARTIST = 1;
  }
  optional Type type = 1;
  optional string identifier = 2;
  optional string date = 3;

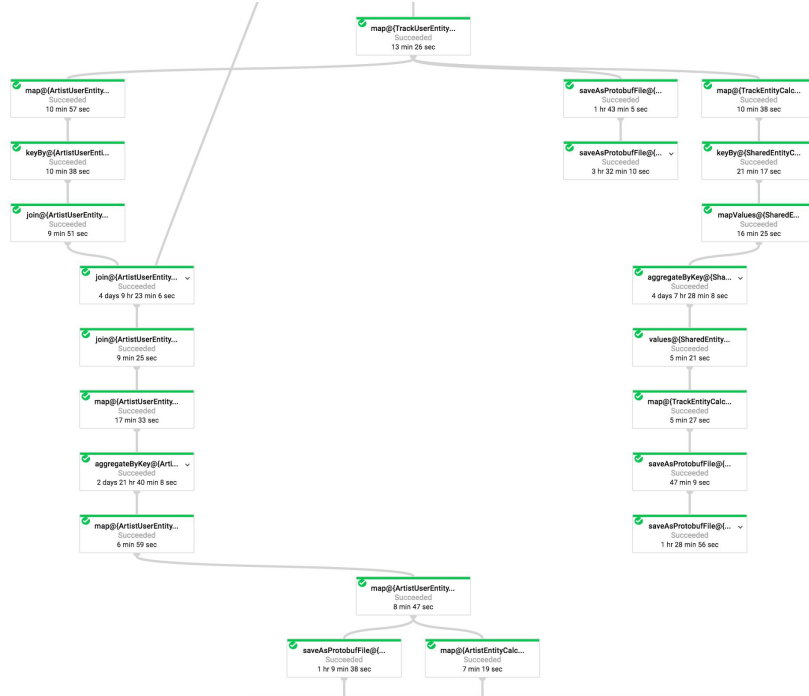
  // Geo and Listener type are part of the key
  optional GeographyInfoPB location = 4;
  optional ListenerTypePB listenerType = 5;
}

// Next ID: 3
message AggregateListenersByGeoPB {
  optional AggregateKeyPB key = 1;
  optional DataByTimeframePB listeners = 2;
}
```

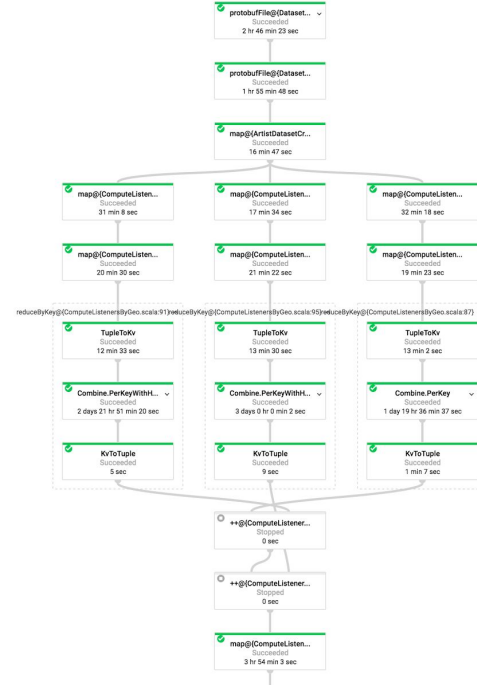
Architecture



Entities



Datasets



Entities Layer



+



of Skips, Saves, Streams

Key

Gender

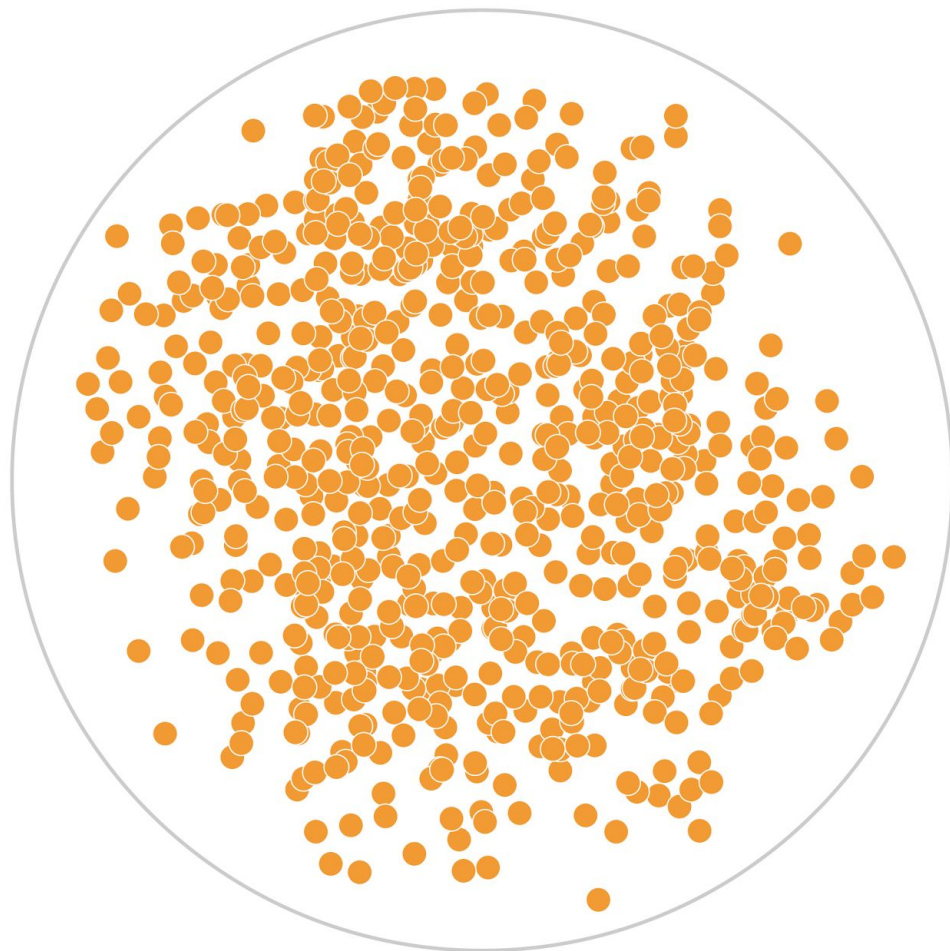
Age

Time

Location

Listener Type

Play Source



Gender

Male

Age

19-24

Time

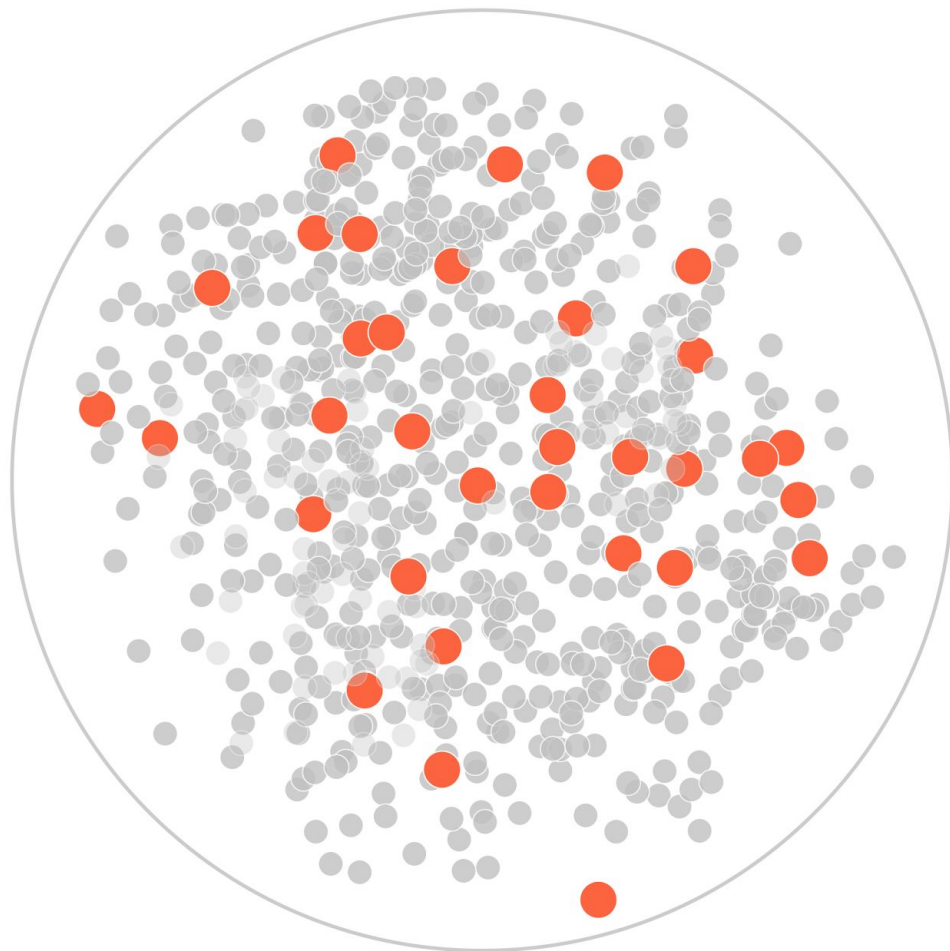
Location

Brazil

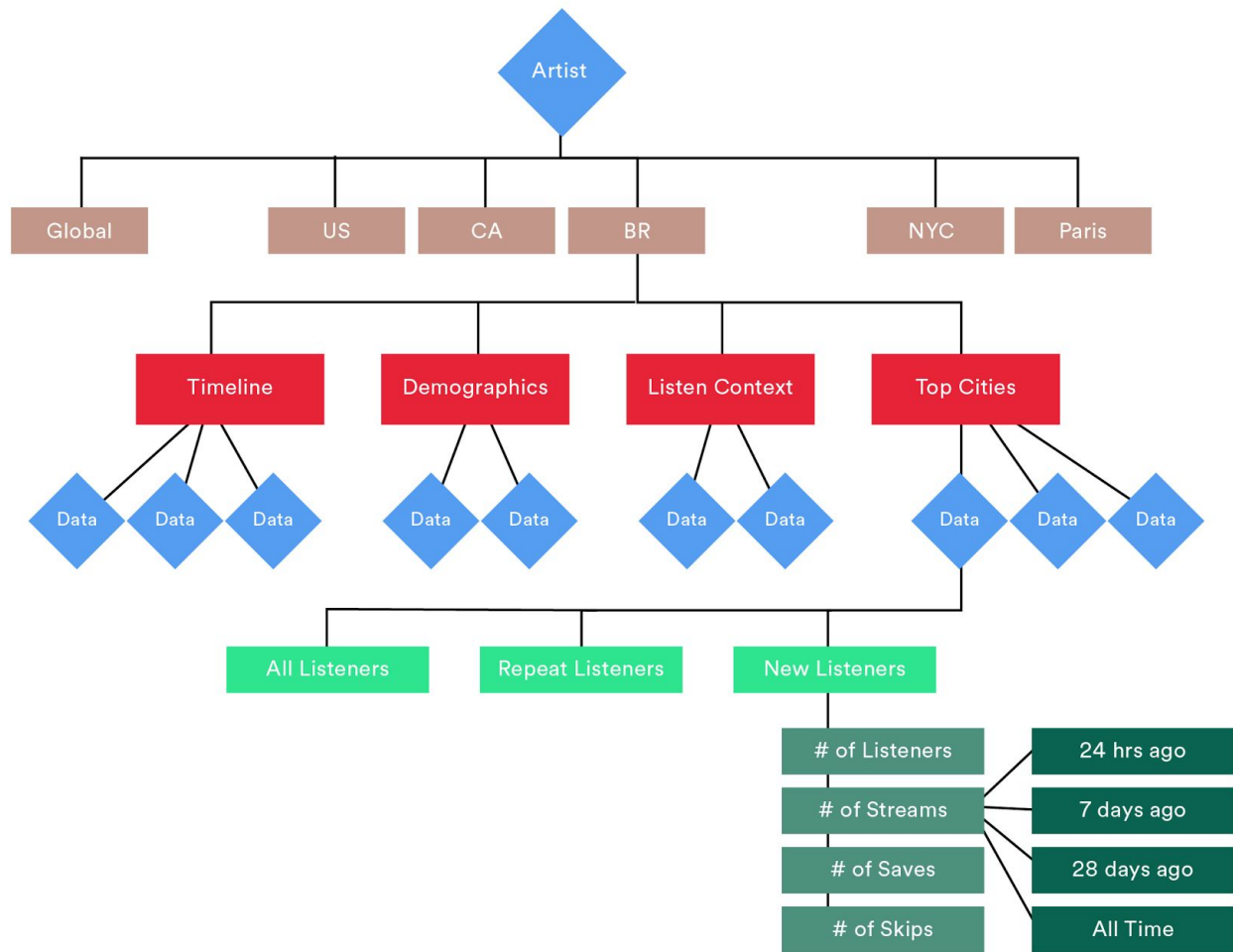
Listener Type

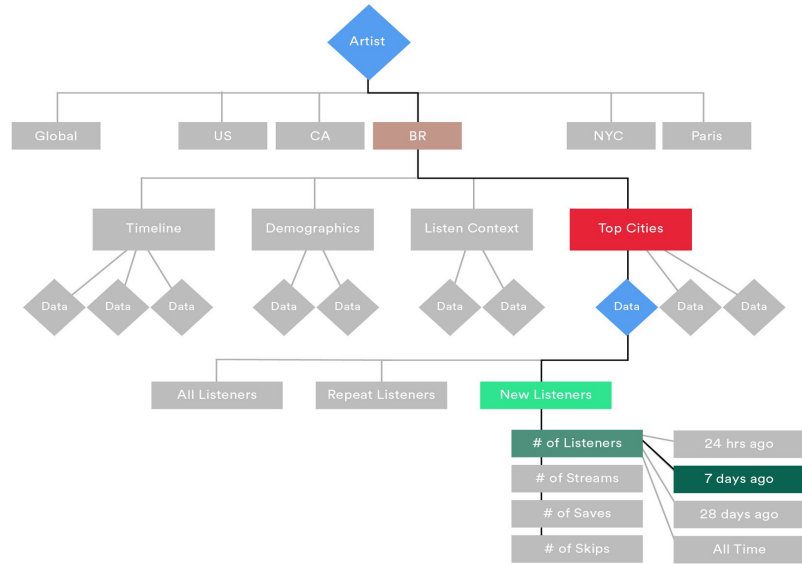
Fans

Play Source



Dataset Layer



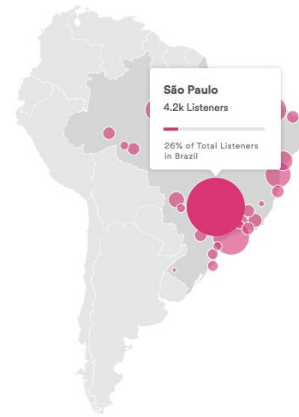


Where they are

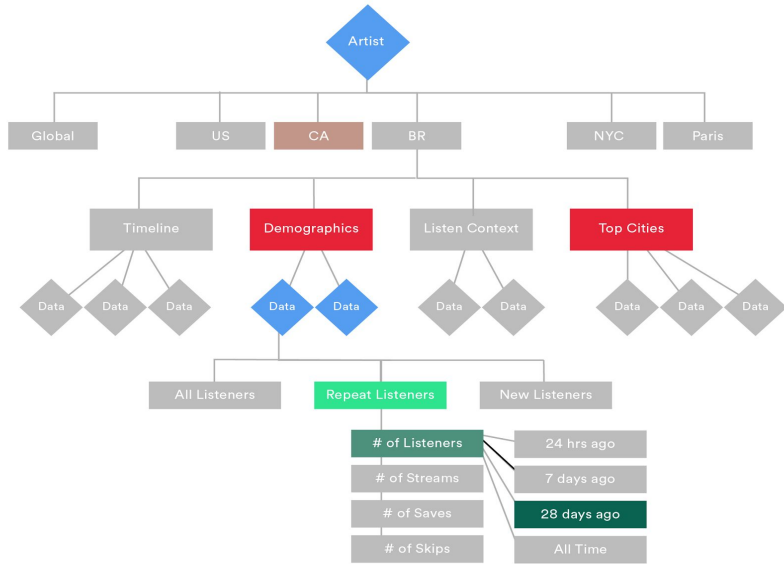
LAST 7 DAYS • BRAZIL

24,573 Listeners

4% of Total Listeners

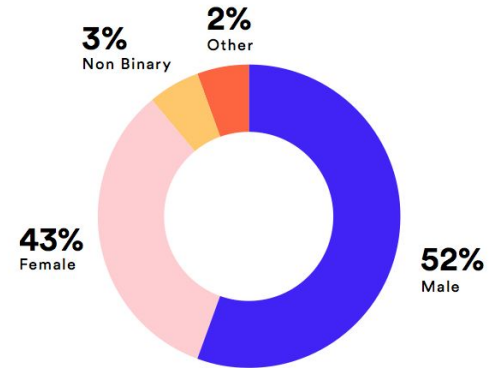


São Paulo	4.2k
Rio de Janeiro	3.1k
Salvador	2.5k
Brasília	2.4k
Fortaleza	1.8k
Belo Horizonte	1.7k
Manaus	1.3k
Curitiba	1.2k
Recife	1.1k
Porto Alegre	827



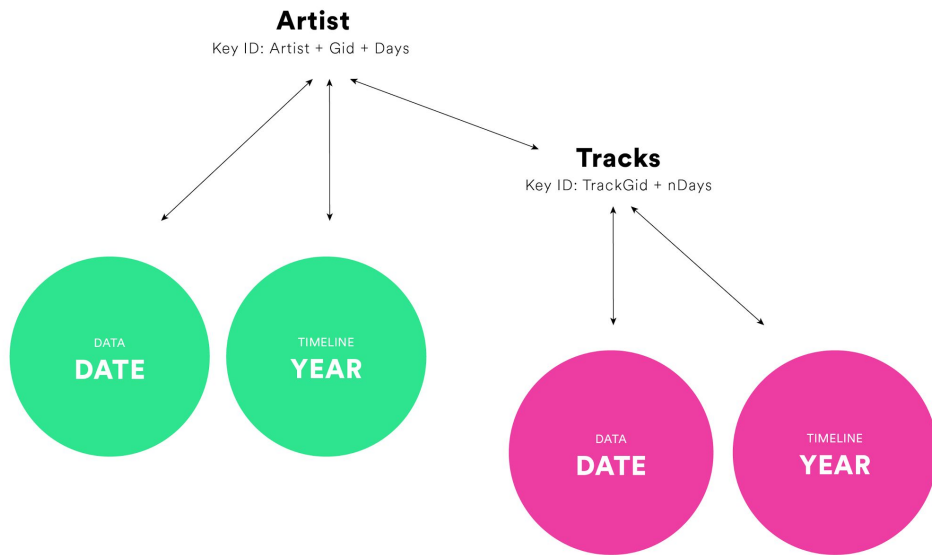
Who they are

LAST 28 DAYS • CANADA



Export Layer

Datastore Storage Model



Conclusions

- Unified Pipeline Architecture allows you to parallelize your processing while consolidating the logic
 - No redundant computations
- Minimize Number of joins
- Read in sources only once and decorate them once for further processing
- Think About the Future!



Questions?

- Come to my office hours!
- My colleague: Deepti Deshpande will also be there to help answer questions