



Universitat  
Pompeu Fabra  
*Barcelona*

**MTG**  
Music Technology  
Group

---

# Knowledge Extraction and Feature Learning for Music Recommendation in the Long Tail

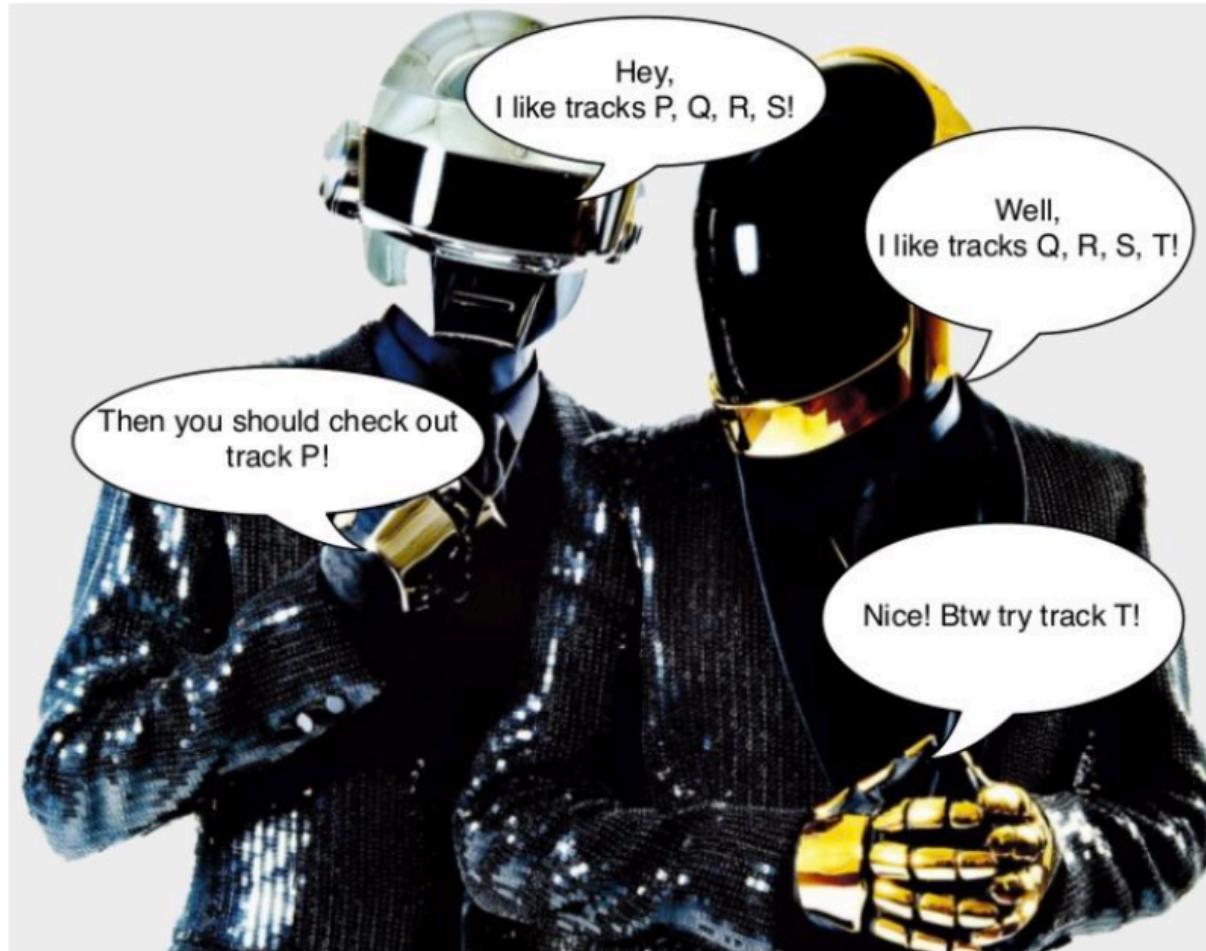
Sergio Oramas

Collaborators: Oriol Nieto, Vito Claudio Ostuni, Tommaso Di Noia, Mohamed Sordo, Xavier Serra



Politecnico di Bari

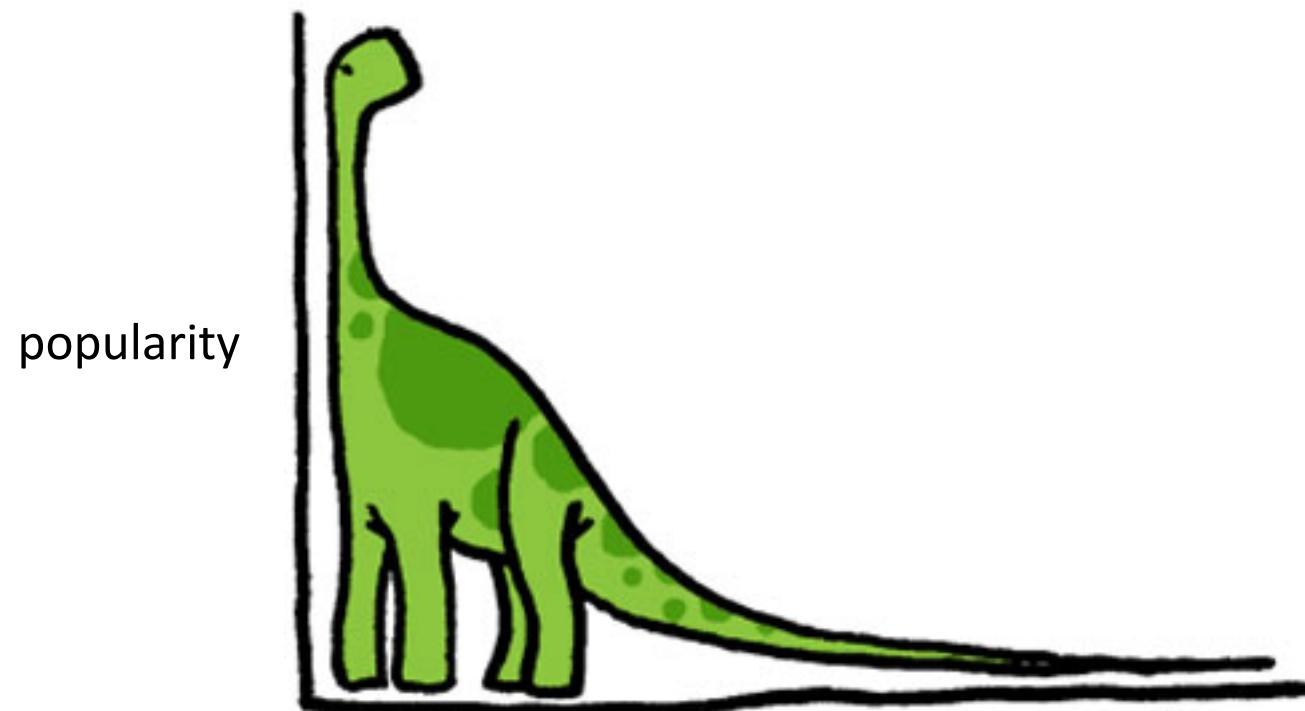
# Collaborative Filtering



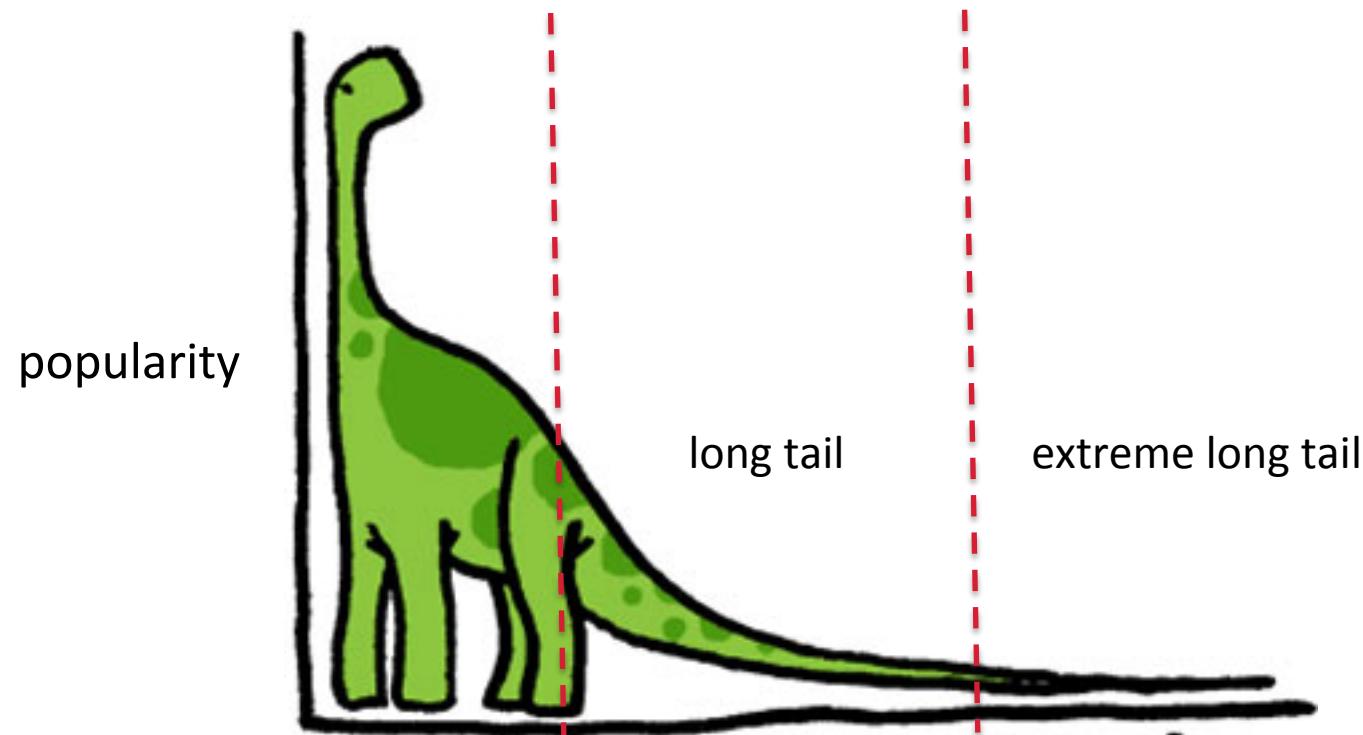
# Collaborative Filtering



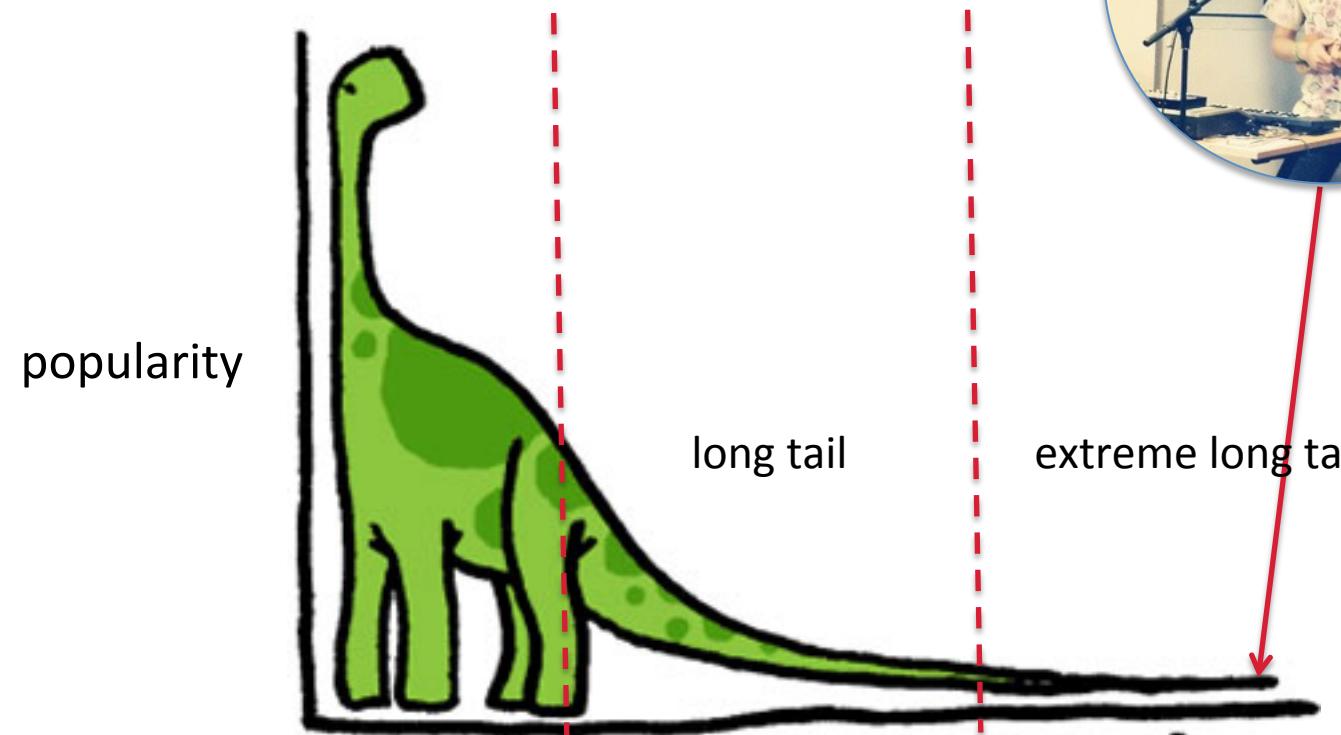
# Long Tail



# Long Tail



# Long Tail



---

# Cold-start



# Cold-start Problem: New Music Releases

- Online Streaming Services daily ingest new releases from:
  - Existing Artists
  - Novel Artists



# Cold-start Problem: Ingest Items in Catalog

- The Pandora case:
  - New on-demand service
  - From millions to tens of millions



# Content features

- **Tags**
    - Widely used
    - Need a curation process (experts or a crowd)
    - Rich contextual info: genres, places, relations, dates



# Content features

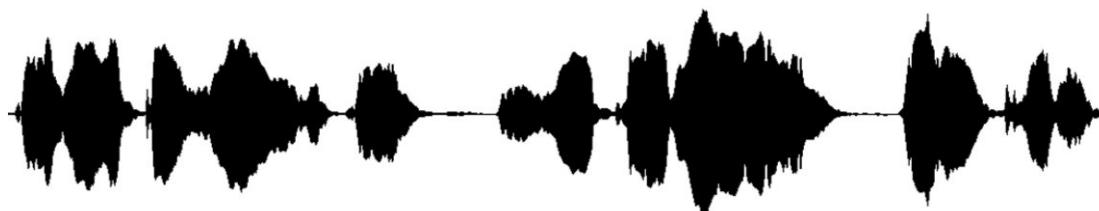
- Artist biographies or press releases
  - Barely used
  - Often self-made
  - Rich contextual info: genres, places, relations, dates
  - Noisy

The screenshot shows a dark-themed website interface. At the top, there is a navigation bar with a play button icon, a double-back arrow, a single-back arrow, a single-forward arrow, and a double-forward arrow. To the right of the arrows is the text "last.fm". Below the navigation bar, the word "Biography" is centered. The main content area contains a detailed biography of the band Wilco. To the right of the biography, there are several sections with headings and associated information. The sections include:

- YEARS ACTIVE**: 1994 – present (23 years)
- FOUNDED IN**: Chicago, Cook County, Illinois, United States
- MEMBERS**: Bob Egan, Glenn Kotche (2000 – present), Jay Bennett (1995 – 2001), Jeff Tweedy

# Content features

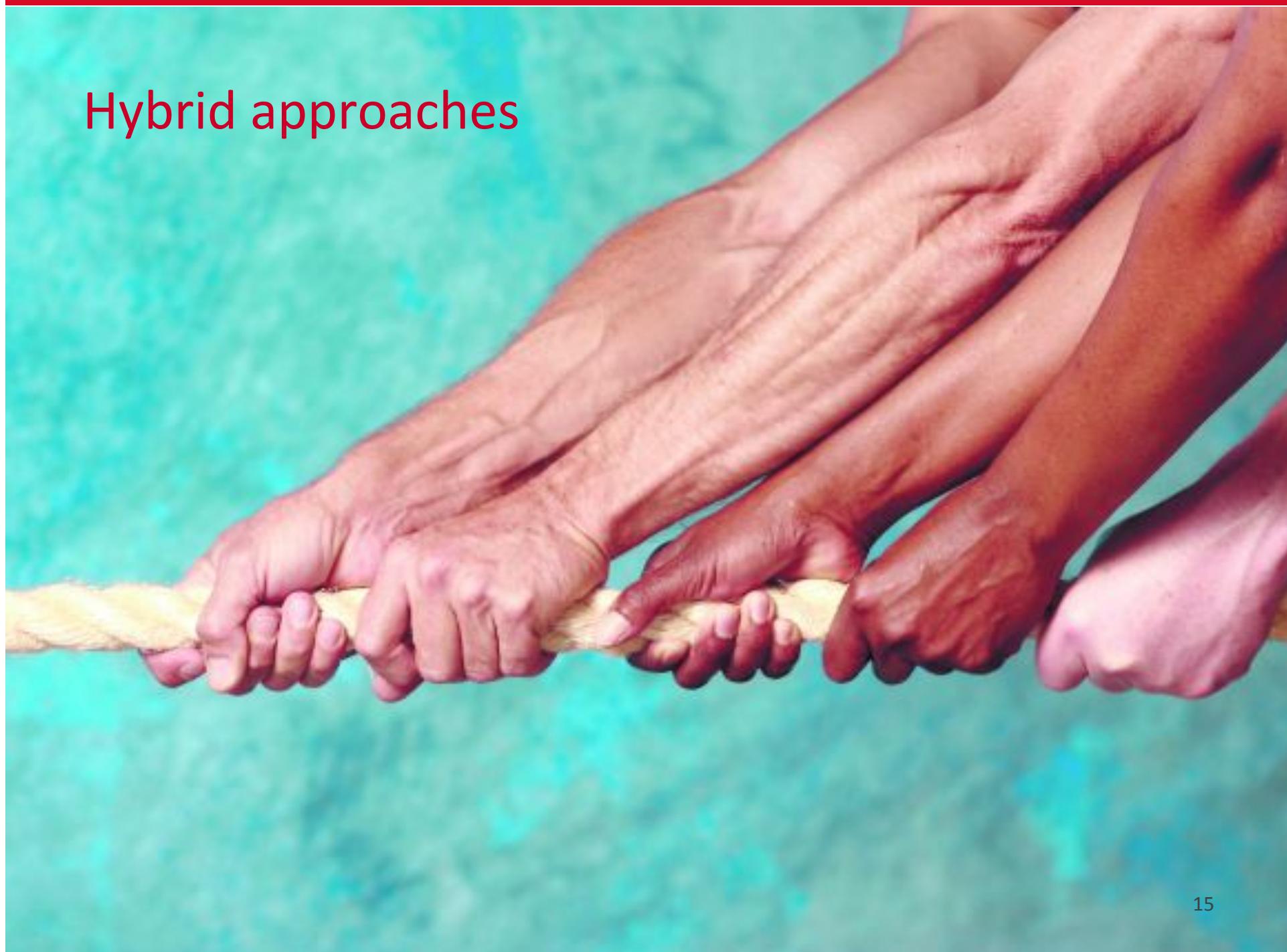
- **Audio**
  - Slightly used
  - Always available
  - Rich musical info: genres, instruments, rhythms, timber
  - Semantic gap



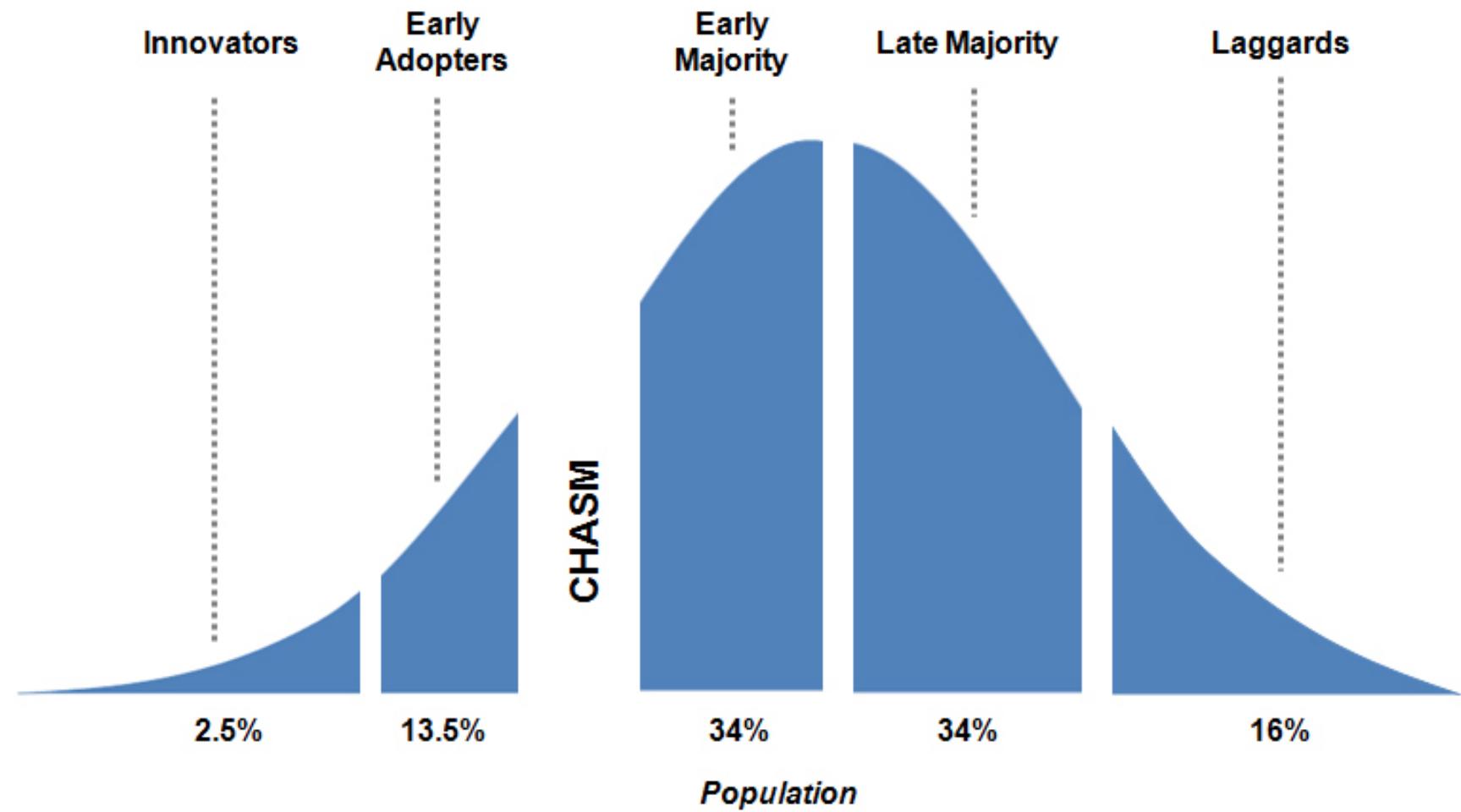
# Content-based vs. Collaborative Filtering



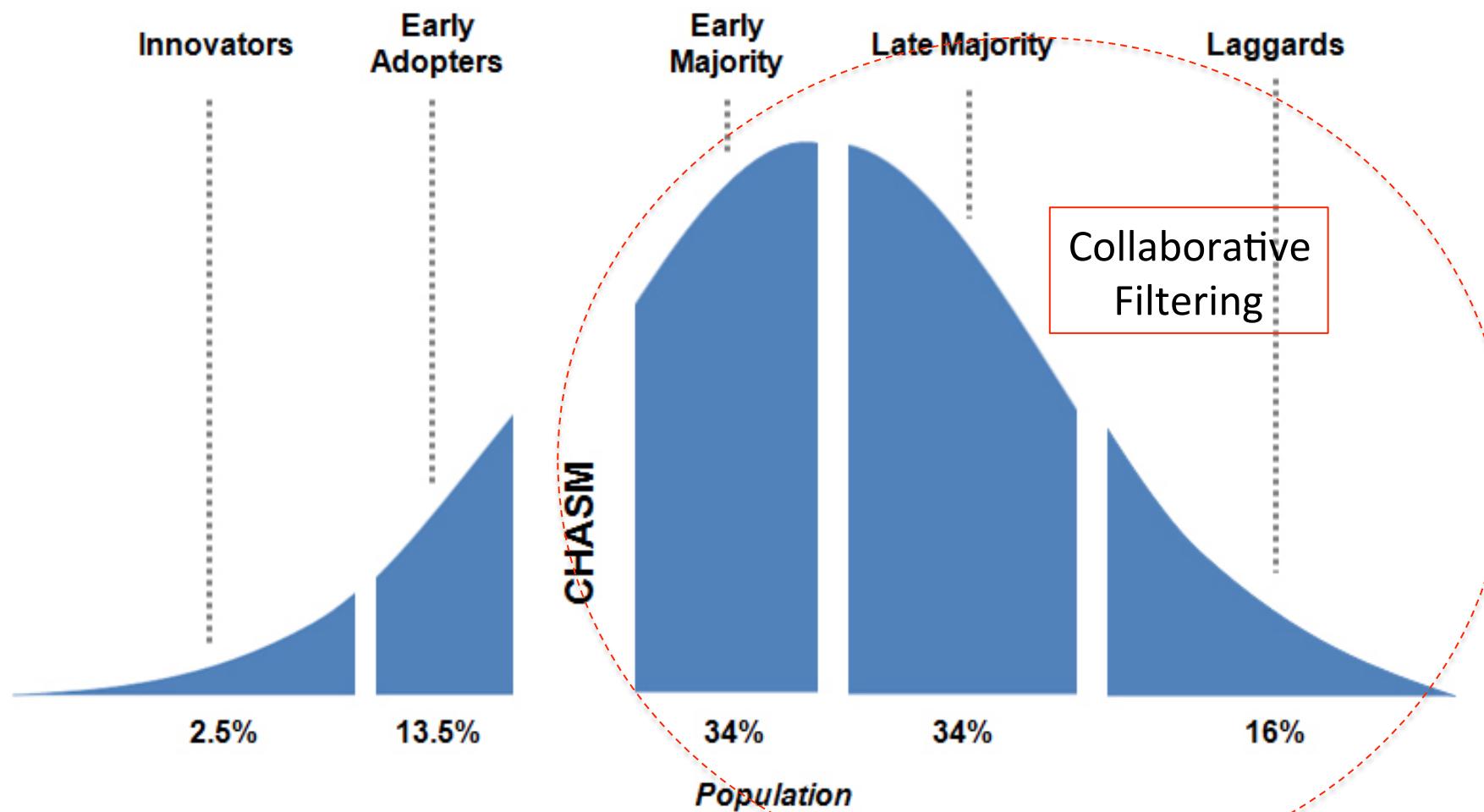
## Hybrid approaches



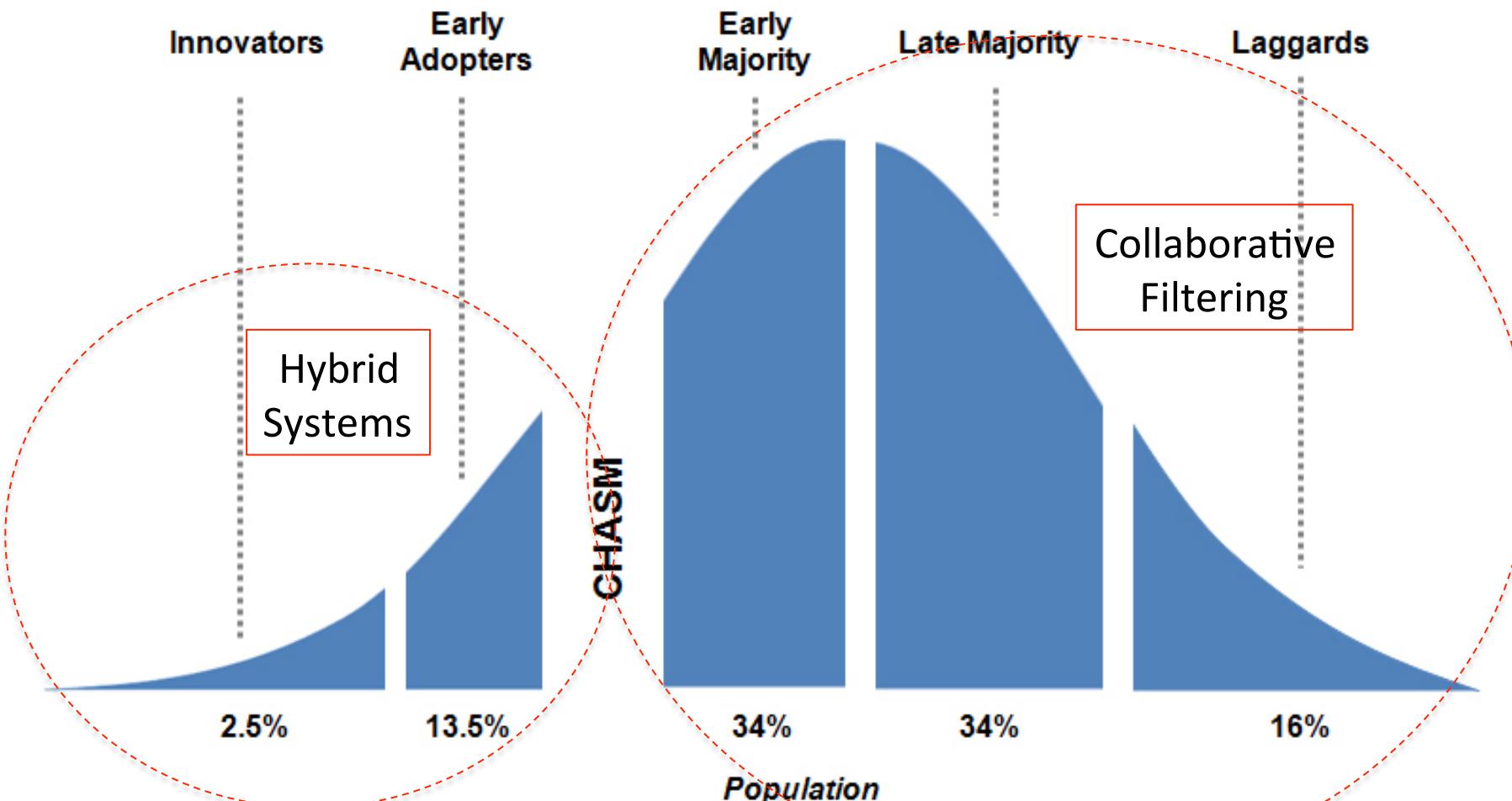
# Users



# Users



# Users



# Knowledge Extraction



Universitat  
Pompeu Fabra  
*Barcelona*

**MTG**  
Music Technology  
Group

# Structured Data Sources

- Structured information about music is incomplete
  - Only popular artists and western music
  - Only editorial and some biographical information



# Unstructured Data Sources

- Huge amount of music information remains implicit in unstructured texts
  - Artists biographies
  - Articles
  - Reviews
  - Web pages
  - User's posts

The screenshot shows a complex interface with several open tabs and windows. The main visible window is from Songfacts.com, displaying details about the song 'Sweet Freedom' by Joe Crepusculo. Below it, a Facebook profile for 'Joe Crepusculo' is visible, showing a profile picture and some status updates. At the bottom, a LinkedIn post is partially visible, dedicated to Joe Crepusculo.

The screenshot displays two Wikipedia pages side-by-side. The left page is for the Nigerian singer Nneka, featuring a biography, a photo, and a 'Factbox' section. The right page is for Karakudi Mani, an Indian percussionist, also with a biography, photo, and 'Factbox'. Both pages include links to other Wikipedia articles and external resources.



# Entity Linking

Identify and link entity mentions in text to a Knowledge Base



Babelfy



# Entity Linking

Songs in Sky Blue Sky were written by Wilco frontman , Jeff Tweedy.

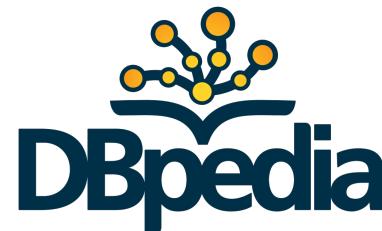
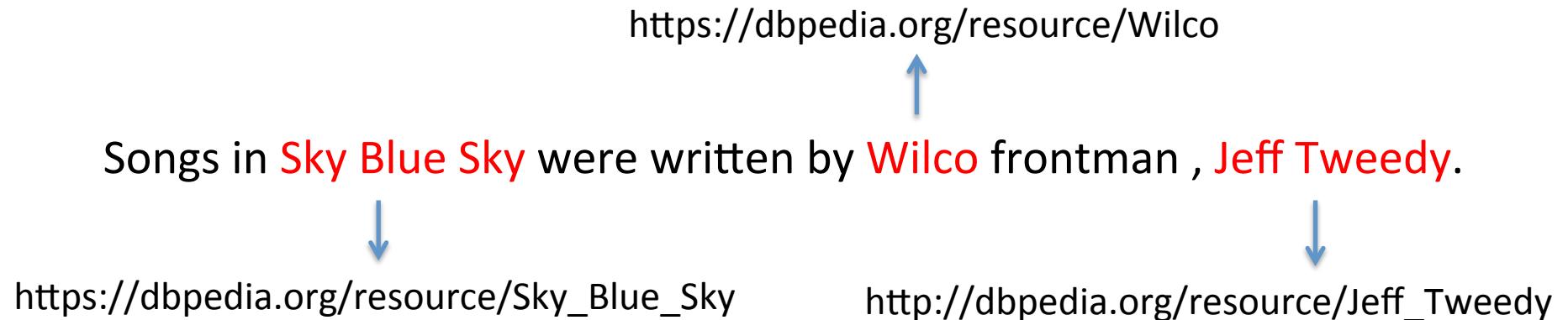


# Entity Linking

Songs in **Sky Blue Sky** were written by **Wilco** frontman , **Jeff Tweedy**.

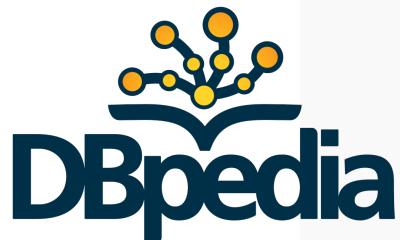


# Entity Linking



# Semantic Information

- Wilco



dbo:genre

- dbr:Art\_rock
- dbr:Experimental\_rock
- dbr:Alternative\_country
- dbr:Alternative\_rock
- dbr:Indie\_rock

dbo:hometown

- dbr:United\_States
- dbr:Chicago
- dbr:Illinois

dbo:recordLabel

- dbr:Reprise\_Records
- dbr:Nonesuch\_Records
- dbr:Dbpm\_records

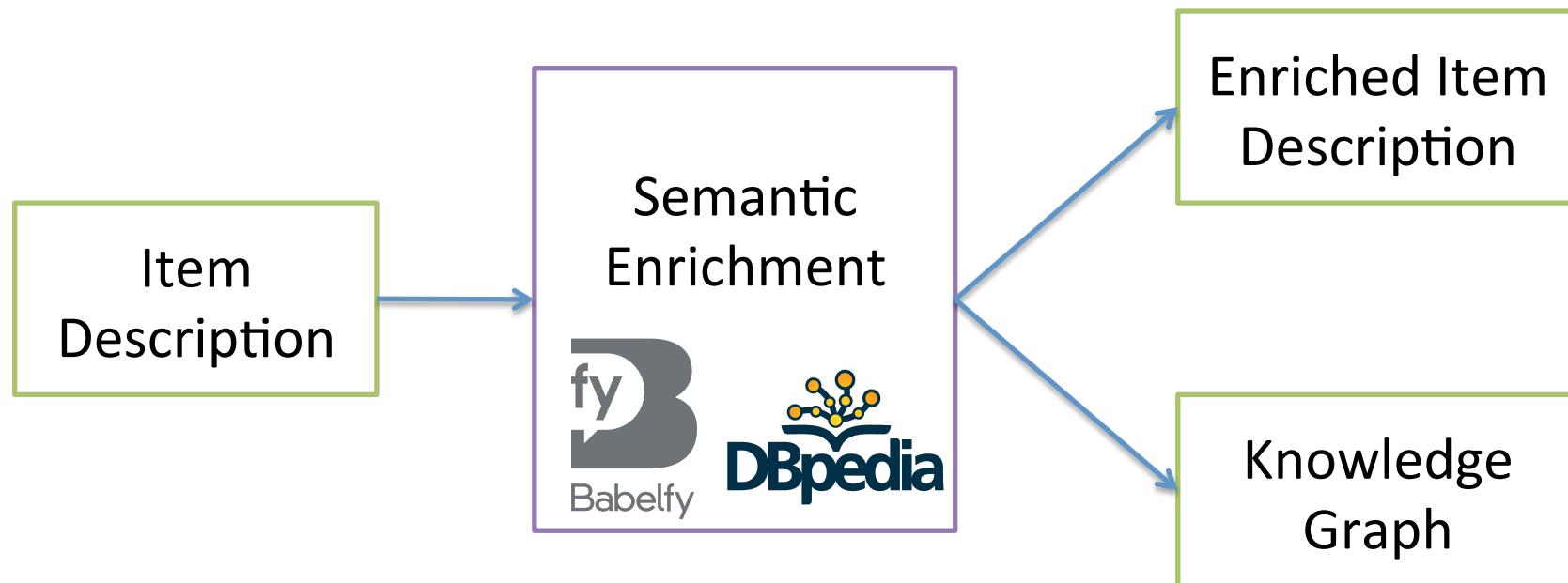


Universitat  
Pompeu Fabra  
Barcelona

MTG

Music Technology  
Group

# Semantic Enrichment via Entity Linking



# Enriched Item Description

Dies ist ein Blindtext. An ihm lässt sich vieles über die Schrift ablesen, in der er gesetzt ist. Auf den ersten Blick wird der Grauwert der Schriftfläche sichtbar. Dann kann man prüfen, wie gut die Schrift zu lesen ist und wie sie auf den Leser wirkt.

Dies ist ein Blindtext. An ihm lässt sich vieles über die Schrift ablesen, in der er gesetzt ist. Auf den ersten Blick wird der Grauwert der Schriftfläche sichtbar. Dann kann man prüfen, wie gut die Schrift zu lesen ist und wie sie auf den Leser wirkt.



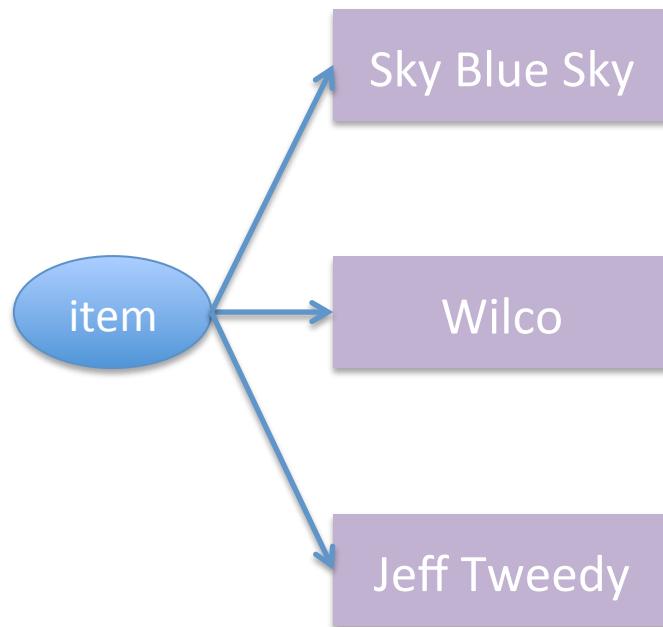
dbr:Art\_rock dbr:Experimental\_rock  
d b r :A l t e r n a t i v e \_c o u n t r y  
dbr:Alternative\_rock dbr:Indie\_rock  
dbr:United\_States dbr:Chicago  
dbr:Illinois dbr:Reprise\_Records  
d b r :N o n e s u c h \_R e c o r d s  
dbr:Dbpm\_records

Biography text

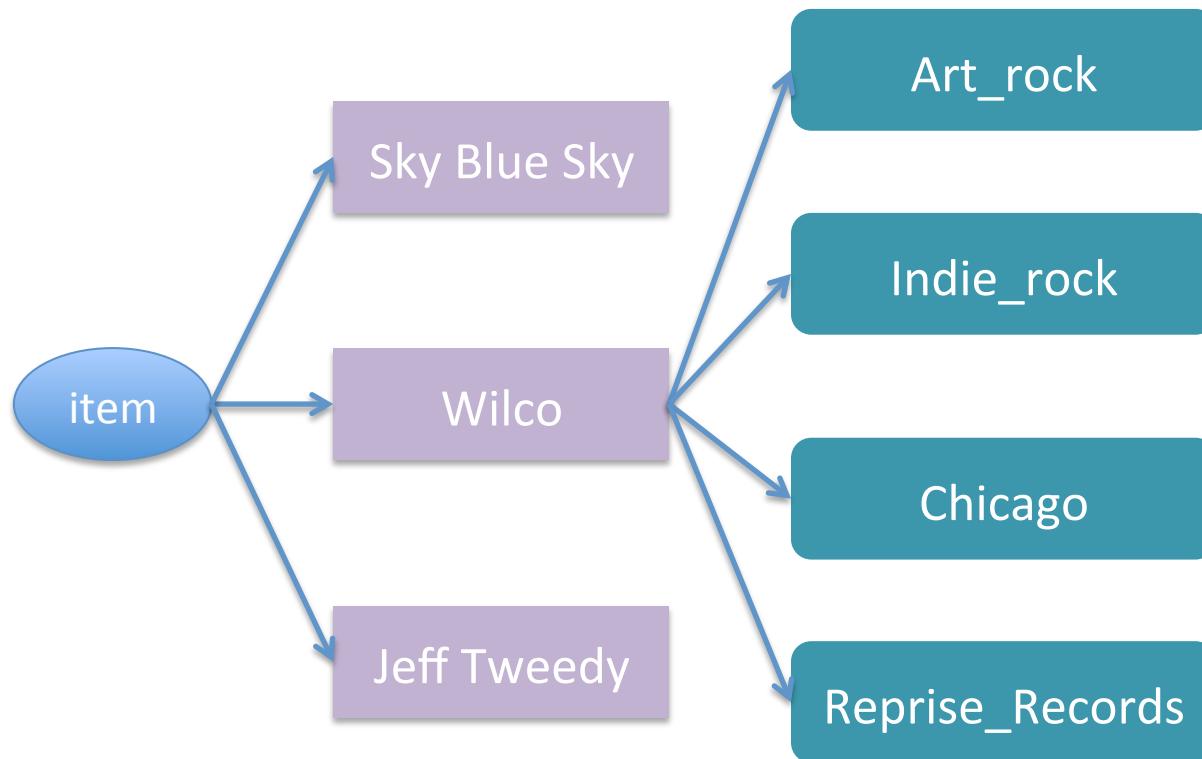
VSM  
tf-idf

Semantic data

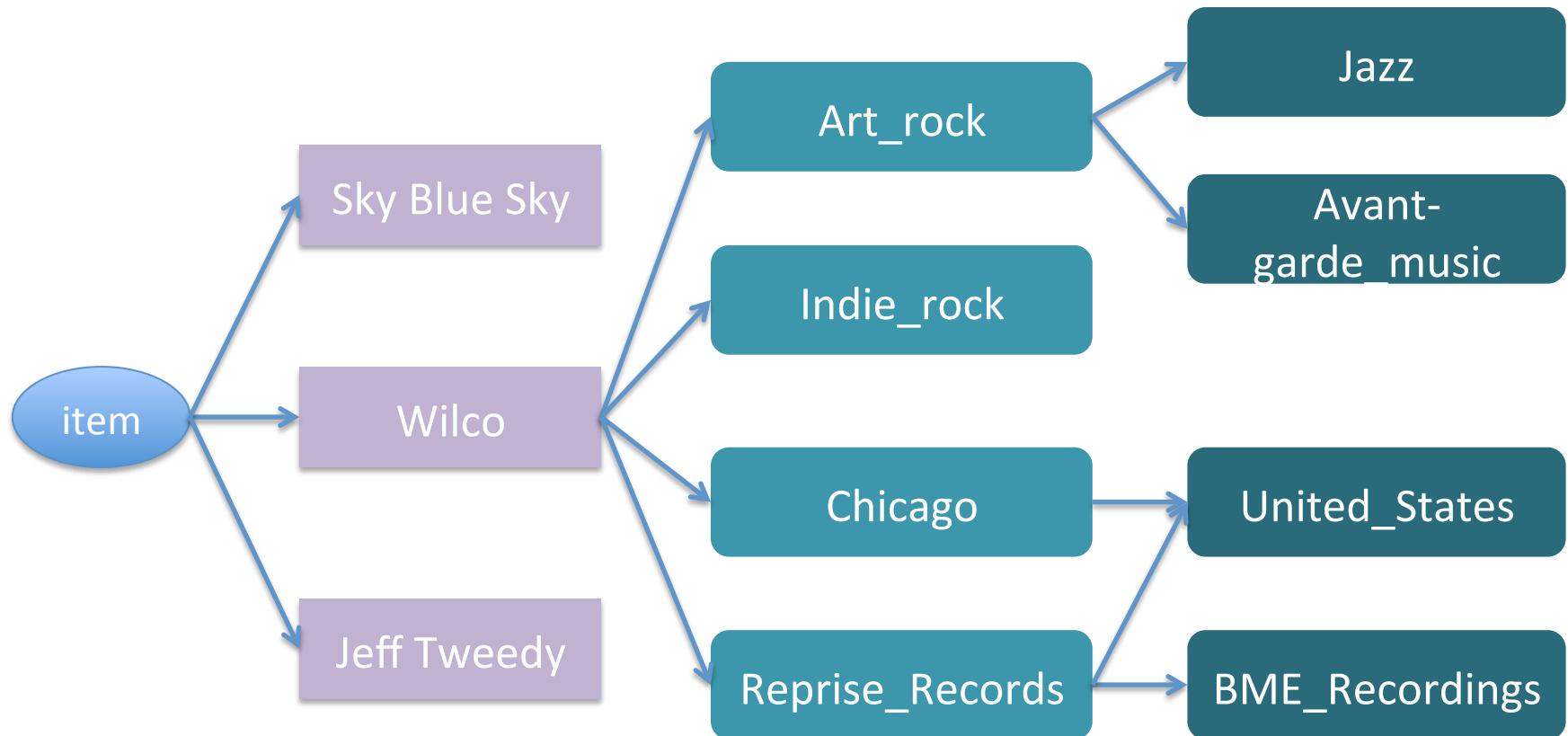
# Knowledge Graph



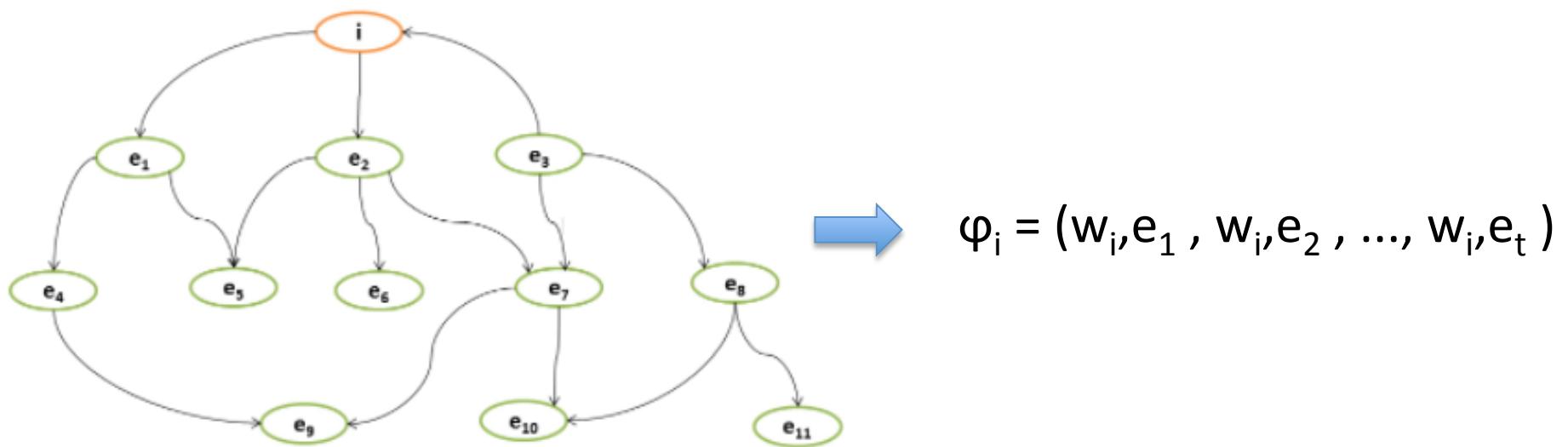
# Knowledge Graph



# Knowledge Graph



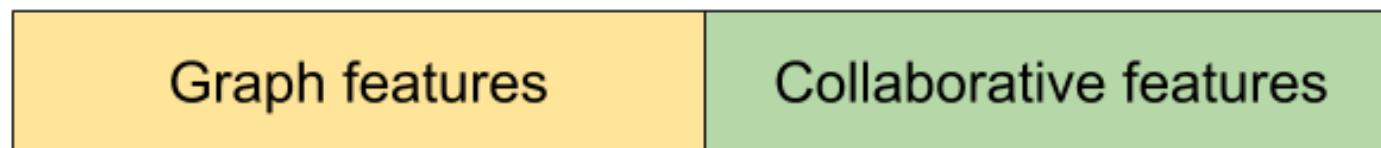
# Graph Embedding



# Long Tail Recommendations

Oramas S., Ostuni V.C., Di Noia T. et al. (2016). **Sound and Music Recommendation with Knowledge Graphs**. ACM-TIST

- Hybrid feature-combination approach
  - Knowledge Graph embedding
  - Collaborative features (user implicit feedback)



# Experiments

- Sounds Recommendation
  - Freesound tags and descriptions
  - Implicit feedback (downloads)
  - 21,552 items and 20,000 users
- Music Recommendation
  - Last.fm tags and Songfacts descriptions
  - Implicit feedback (Last.fm listening habits)
  - 8,640 items and 5,199 users

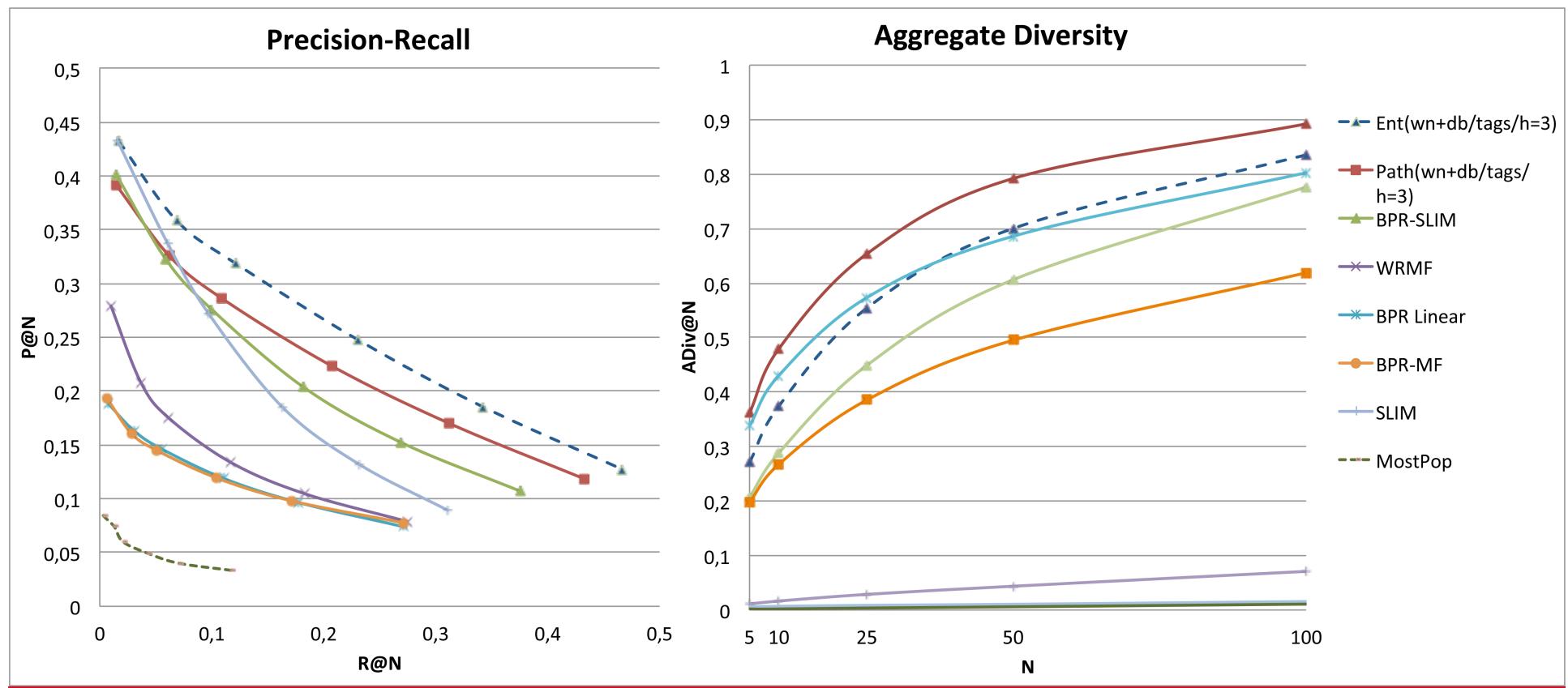
# Results

- Semantic Enrichment in hybrid approach
  - Keeping high accuracy improve novelty and diversity
  - Better explore long tail items

Approach	Accuracy	Novelty	Diversity
Collab			
Collab + Tags			
Collab + Sem			
Sem only			

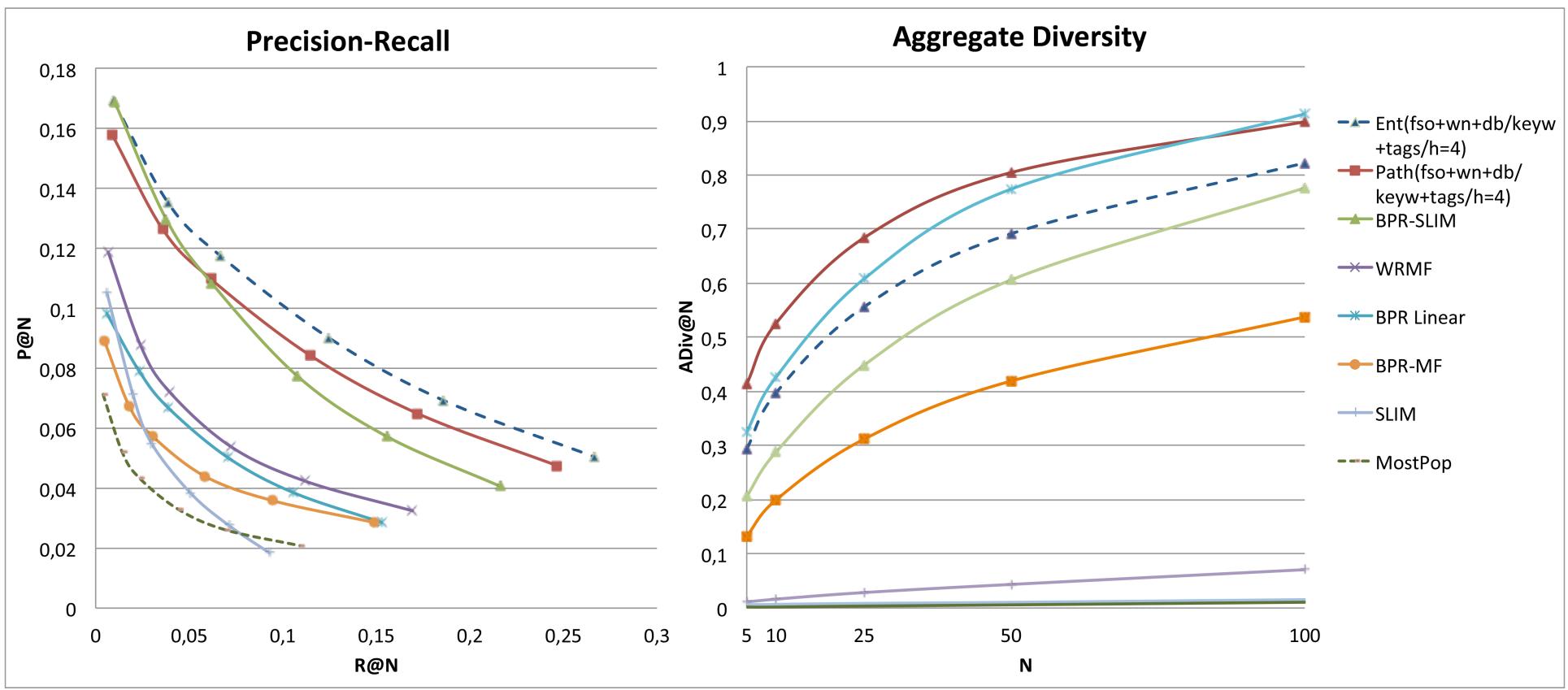
# Results

- Music Recommendation



# Results

- Sound Recommendation



# Feature Learning



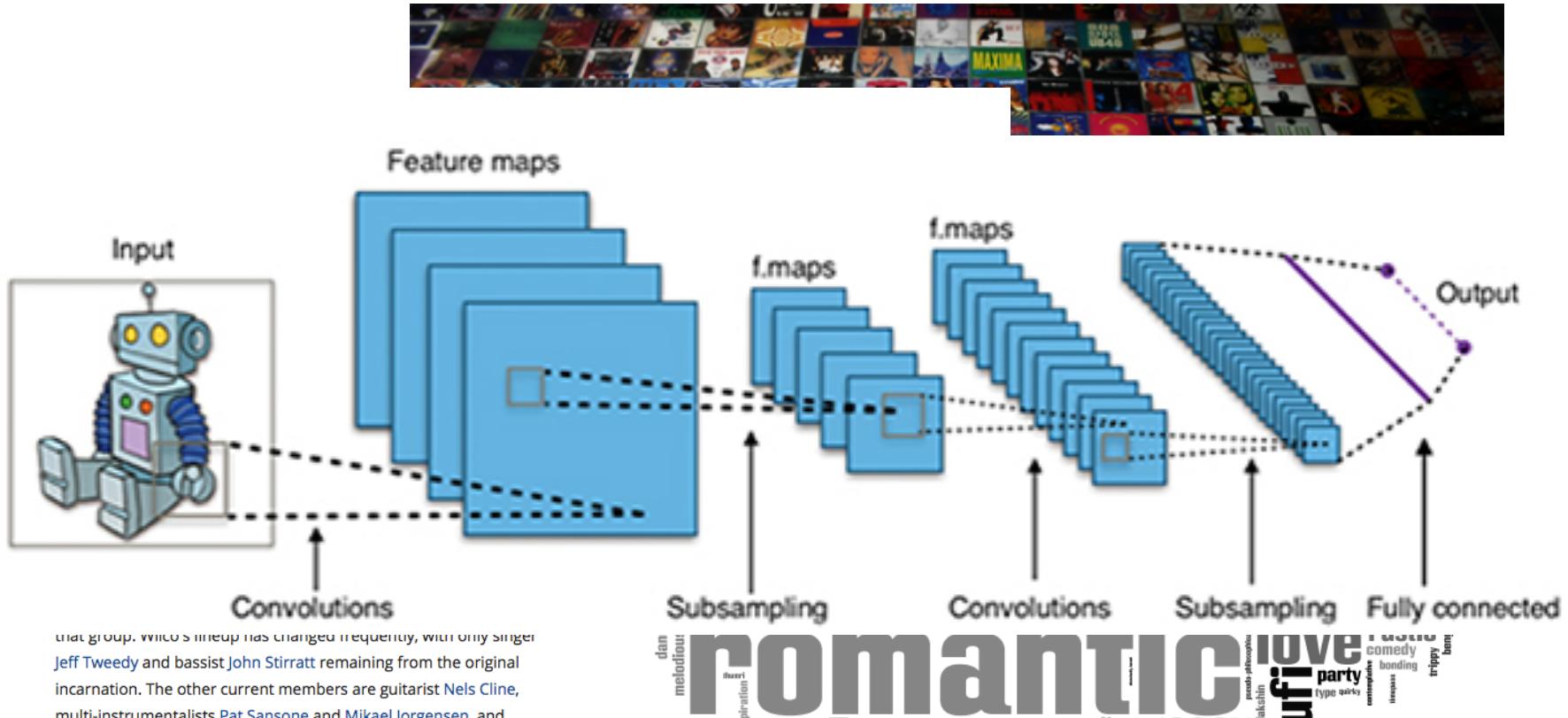
Universitat  
Pompeu Fabra  
*Barcelona*

**MTG**  
Music Technology  
Group

# Multimodal Data



# Multimodal Data



Wilco's lineup has changed frequently, with only singer Jeff Tweedy and bassist John Stiratt remaining from the original incarnation. The other current members are guitarist Nels Cline, multi-instrumentalists Pat Sansone and Mikael Jorgensen, and drummer Glenn Kotche. Wilco has released nine studio albums, a live double album, and three collaborations: two with [Billy Bragg](#), and one with [The Minus 5](#).

# Cold-start Recommendation

Sergio Oramas, Oriol Nieto, Mohamed Sordo, Xavier Serra (2017)  
**A Deep Multimodal Approach for Cold-start Music  
Recommendation.** DLRS Workshop, RecSys 2017.

# Divide & Conquer

## Song features

The Beatles	Love me do
The Beatles	Let it be
The Beatles	A day in the life

# Divide & Conquer

## Song features

The Beatles	Love me do
The Beatles	Let it be
The Beatles	A day in the life

## Artist features

The Beatles

## Track features

Love me do  
Let it be  
A day in the life



# Divide & Conquer



## Song features

The Beatles	Love me do
The Beatles	Let it be
The Beatles	A day in the life



## Track features

Love me do
Let it be
A day in the life

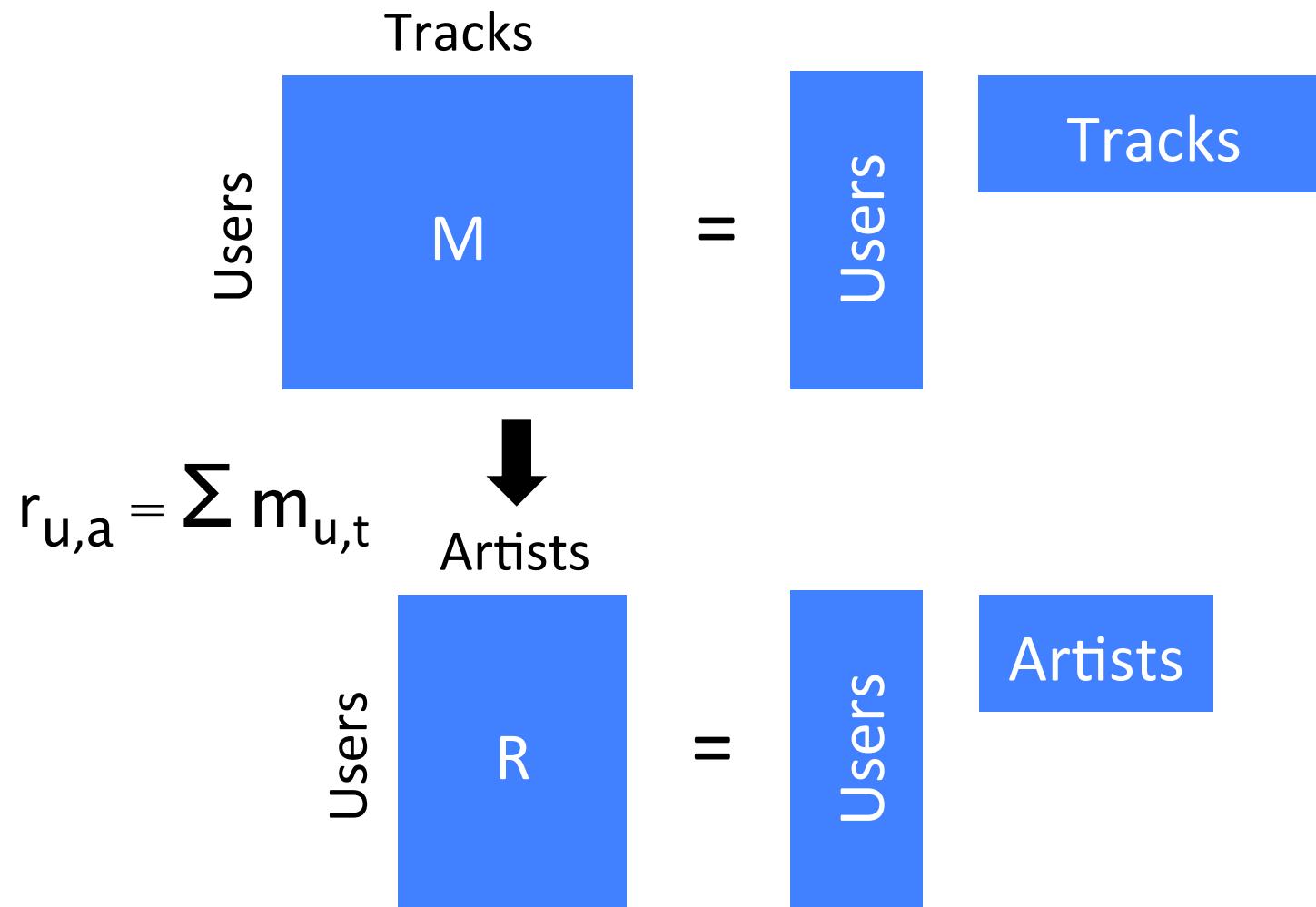
## Artist features

The Beatles
-------------

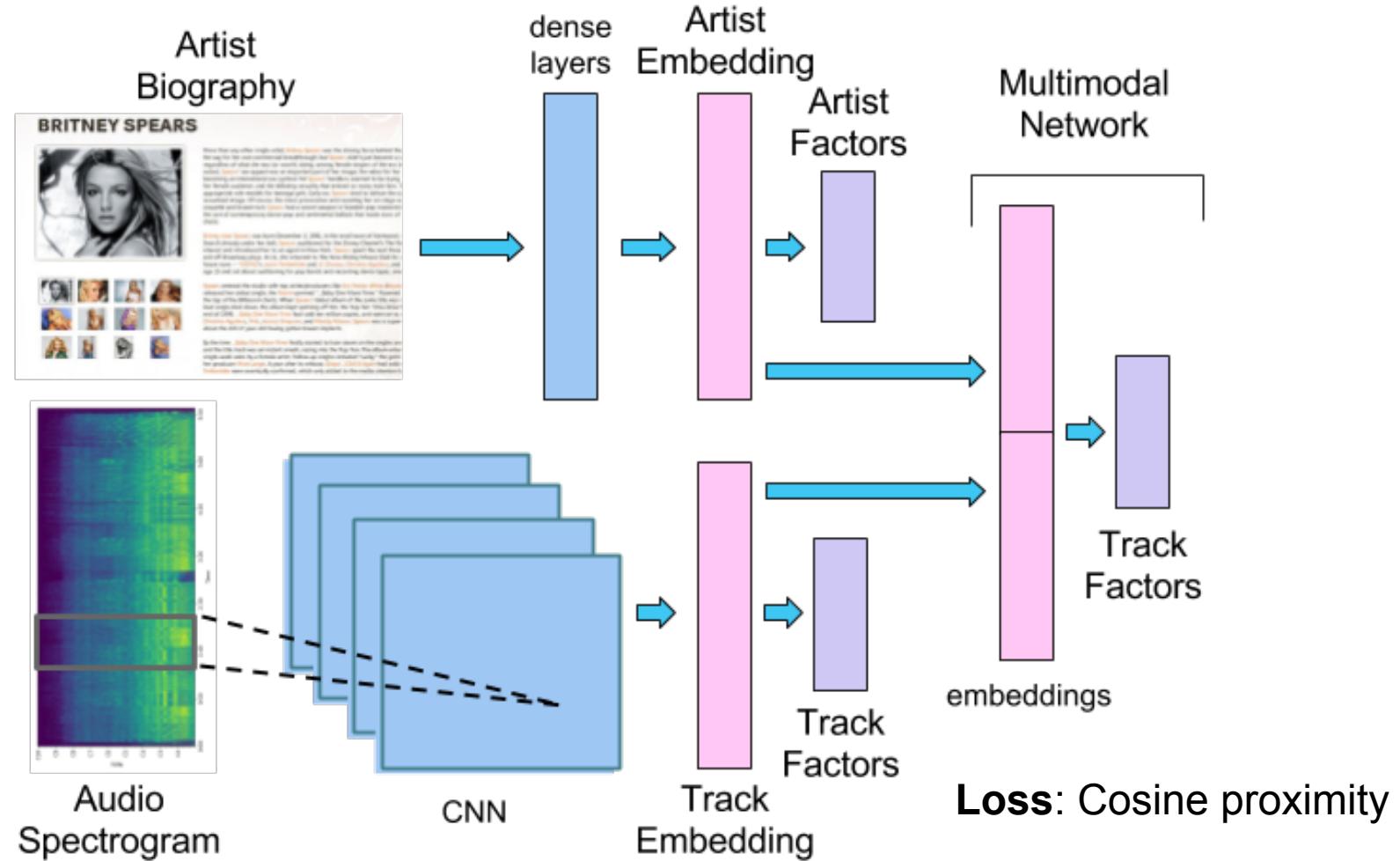
# Recommendation Approach

1. Aggregate feedback data by artist
2. Learn artist feature embeddings from text
3. Learn track feature embeddings from audio
4. Fusion of feature embeddings

# Matrix Factorization (WMF)



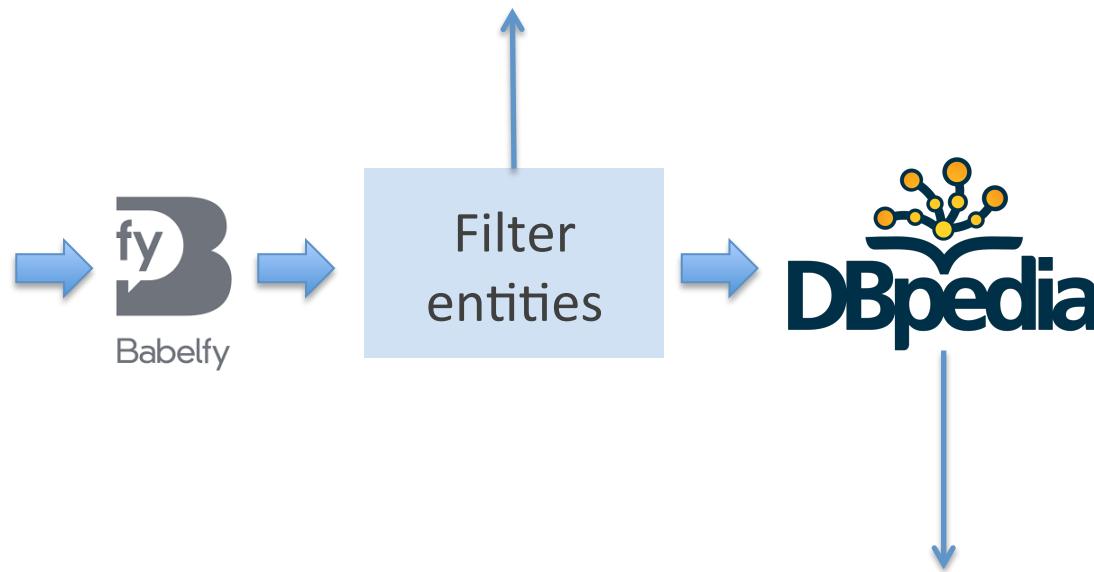
# Recommendation Approach



# Semantic Enrichment

**Types:** MusicalArtist, Band, MusicGenre, MusicalWork, RecordLabel, Instrument, Engineer, and Place

Dies ist ein Blindtext. An ihm lässt sich vieles über die Schrift ablesen, in der er gesetzt ist. Auf den ersten Blick wird der Grauwert der Schriftfläche sichtbar. Dann kann man prüfen, wie gut die Schrift zu lesen ist und wie sie auf den Leser wirkt. Dies ist ein Blindtext. An ihm lässt sich vieles über die Schrift ablesen, in der er gesetzt ist. Auf den ersten Blick wird der Grauwert der Schriftfläche sichtbar. Dann kann man prüfen, wie gut die Schrift zu lesen ist und wie sie auf den Leser wirkt.



**Properties:** homeTown, instrument, genre, associatedBand, writer, producer, recordedIn, etc.

# Artist Text Embeddings: Semantic Enrichment

Dies ist ein Blindtext. An ihm lässt sich vieles über die Schrift ablesen, in der er gesetzt ist. Auf den ersten Blick wird der Grauwert der Schriftfläche sichtbar. Dann kann man prüfen, wie gut die Schrift zu lesen ist und wie sie auf den Leser wirkt. Dies ist ein Blindtext. An ihm lässt sich vieles über die Schrift ablesen, in der er gesetzt ist. Auf den ersten Blick wird der Grauwert der Schriftfläche sichtbar. Dann kann man prüfen, wie gut die Schrift zu lesen ist und wie sie auf den Leser wirkt.

+

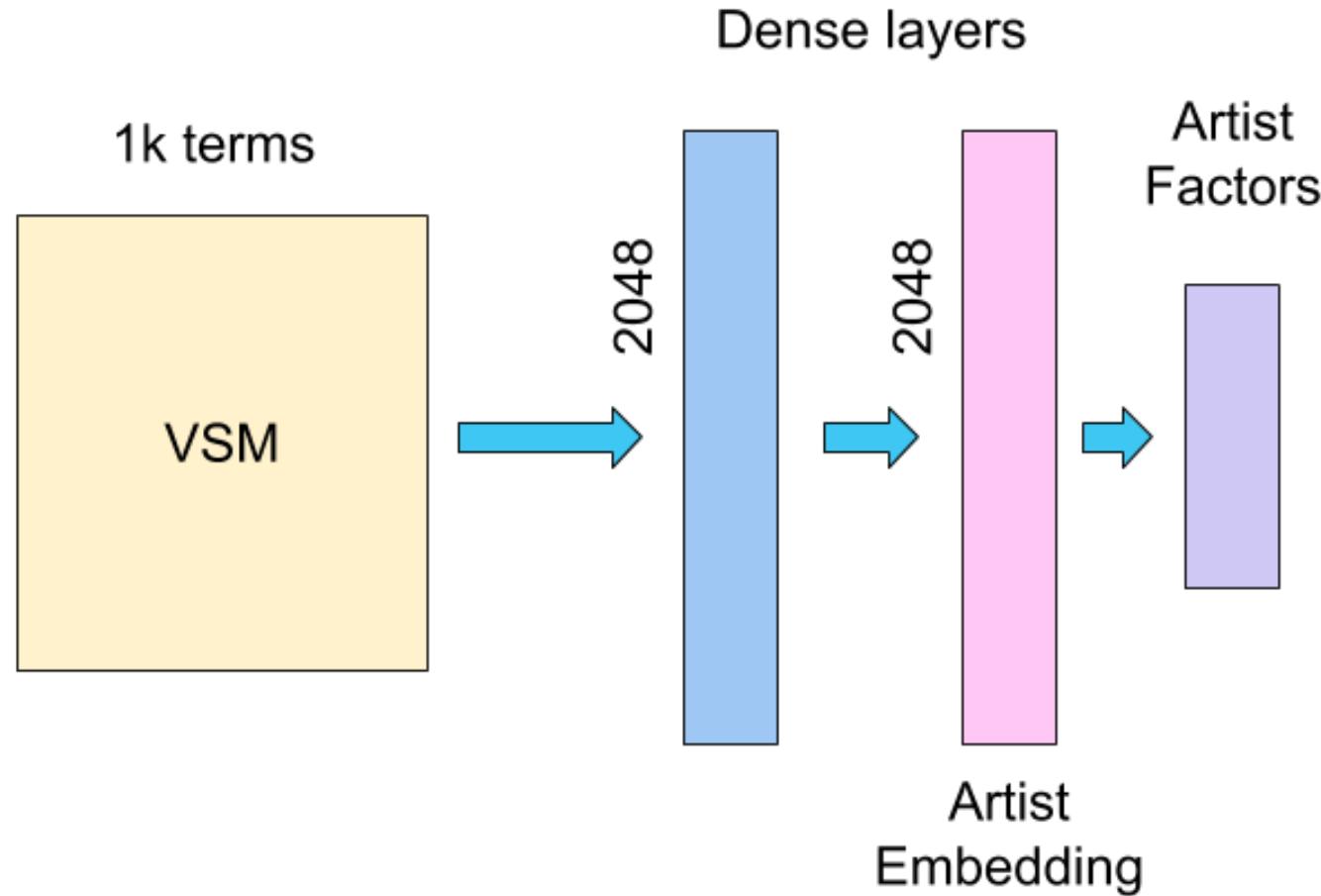
dbr:Art\_rock dbr:Experimental\_rock  
d b r :A l t e r n a t i v e \_ c o u n t r y  
dbr:Alternative\_rock dbr:Indie\_rock  
dbr:United\_States dbr:Chicago  
dbr:Illinois dbr:Reprise\_Records  
d b r :N o n e s u c h \_ R e c o r d s  
dbr:Dbpm\_records

Biography text

VSM  
tf-idf

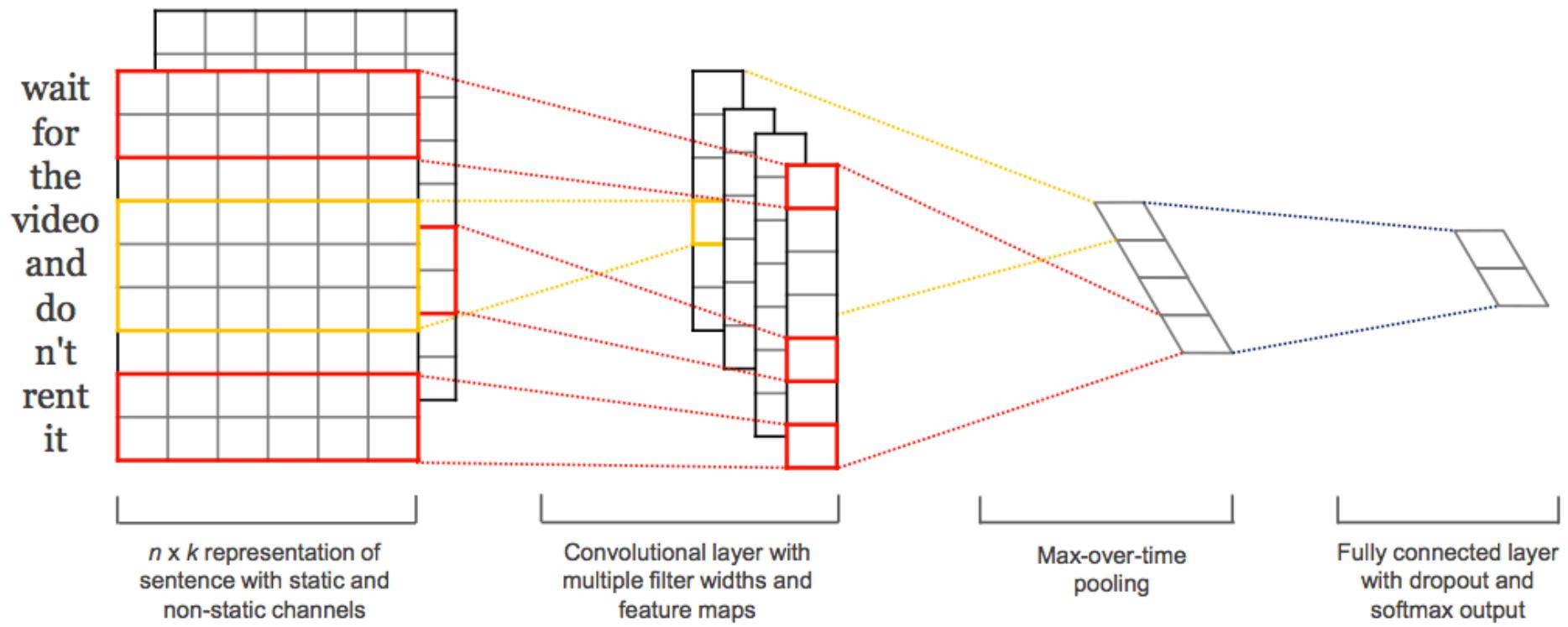
Semantic data

# Artist Text Embeddings: Semantic Enrichment



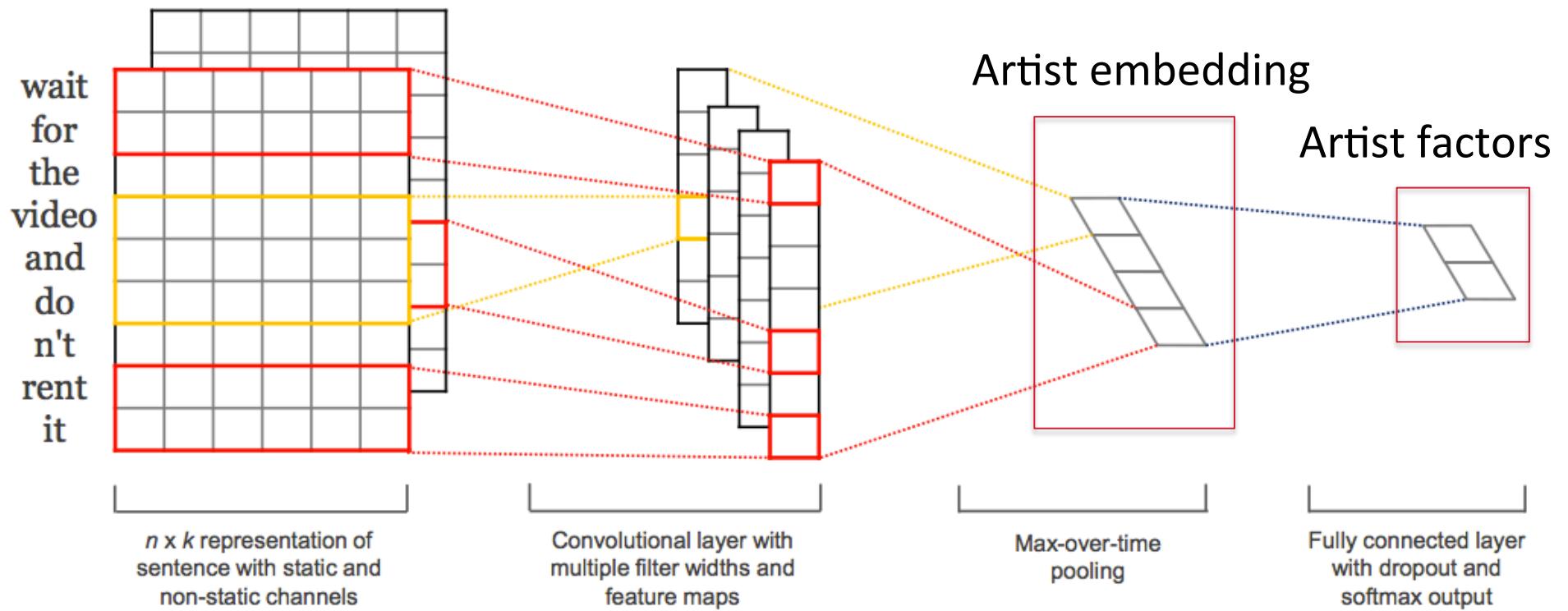
# Artist Text Embeddings: Word Vectors

Kim, Y. (2014). Convolutional neural networks for sentence classification.



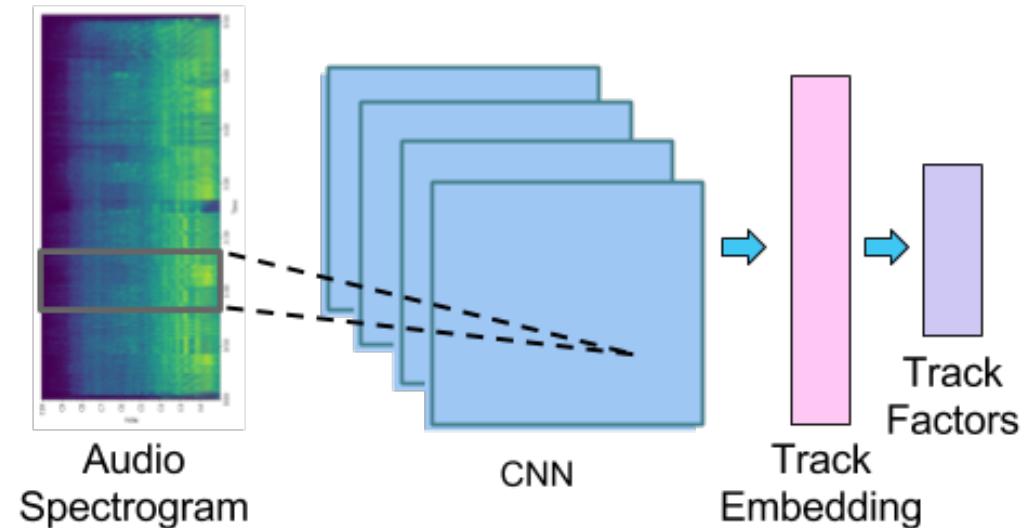
# Artist Text Embeddings: Word Vectors

Kim, Y. (2014). Convolutional neural networks for sentence classification.



# Track Audio Embeddings

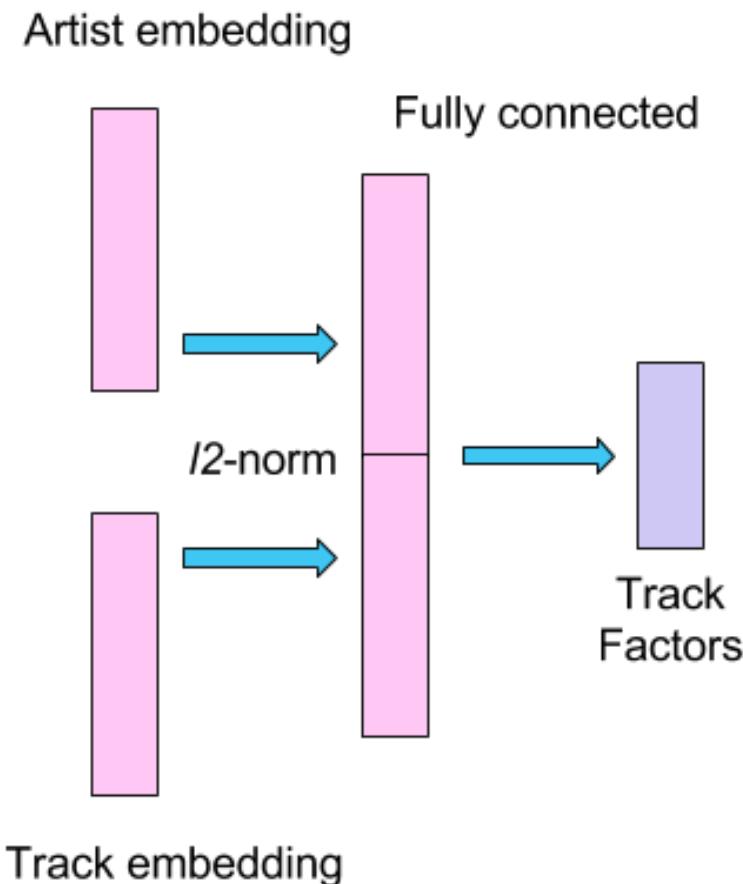
- CQT-spectrograms 96 bins
- 15 sec. patches
- 4 convolutional layers
- Time domain filters
- No dense layers
- Dropout 0.5
- Adam optimizer



Van den Oord et al. (2013) Deep content-based music recommendation

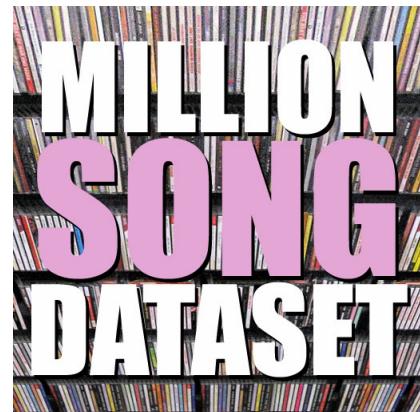
# Multimodal Fusion

- $\ell_2$ -norm
- Concatenate embeddings
- Dropout 0.7



# Dataset

- 328,821 tracks
- 24,043 artists
- ~1M users



Taste Profile

+ last.fm

Artists biographies and tags



Universitat  
Pompeu Fabra  
Barcelona

MTG

Music Technology  
Group

# Experiments

- Artist Recommendation
  - Text-based approaches
- Song Recommendation
  - Audio-based and multimodal approaches
- Splits: 80 train - 10 validation – 10 test
- **Different artists in each subset**
- Evaluation metric: MAP@500

# Artist Recommendation: Text approaches

Aproach	Input	Data model	Arch	MAP
A-TEXT	Bio	VSM	FF	0.0161
<b>A-SEM</b>	<b>Sem Bio</b>	<b>VSM</b>	<b>FF</b>	<b>0.0201</b>
A-W2V-GOO	Bio	w2v-pretrain	CNN	0.0119
A-W2V	Bio	w2v-trained	CNN	0.0145
A-TAGS	Tags	VSM	FF	0.0314
TAGS-ITEMKNN	Tags	-	itemKnn	0.0161
TEXT-RF	Bio	VSM	RF	0.0089
RANDOM	-	-	-	0.0014
UPPER-BOUND	-	-	-	0.5528

# Artist Recommendation: Text approaches

Aproach	Input	Data model	Arch	MAP
A-TEXT	Bio	VSM	FF	0.0161
<b>A-SEM</b>	<b>Sem Bio</b>	<b>VSM</b>	<b>FF</b>	<b>0.0201</b>
A-W2V-GOO	Bio	w2v-pretrain	CNN	0.0119
A-W2V	Bio	w2v-trained	CNN	0.0145
A-TAGS	Tags	VSM	FF	0.0314
TAGS-ITEMKNN	Tags	-	itemKnn	0.0161
TEXT-RF	Bio	VSM	RF	0.0089
RANDOM	-	-	-	0.0014
UPPER-BOUND	-	-	-	0.5528

# Artist Recommendation: Tags vs. Bios

Aproach	Input	Data model	Arch	MAP
A-TEXT	Bio	VSM	FF	0.0161
<b>A-SEM</b>	<b>Sem Bio</b>	<b>VSM</b>	<b>FF</b>	<b>0.0201</b>
A-W2V-GOO	Bio	w2v-pretrain	CNN	0.0119
A-W2V	Bio	w2v-trained	CNN	0.0145
A-TAGS	Tags	VSM	FF	0.0314
TAGS-ITEMKNN	Tags	-	itemKnn	0.0161
TEXT-RF	Bio	VSM	RF	0.0089
RANDOM	-	-	-	0.0014
UPPER-BOUND	-	-	-	0.5528

# Artist Recommendation: Competitor

Aproach	Input	Data model	Arch	MAP
A-TEXT	Bio	VSM	FF	0.0161
<b>A-SEM</b>	<b>Sem Bio</b>	<b>VSM</b>	<b>FF</b>	<b>0.0201</b>
A-W2V-GOO	Bio	w2v-pretrain	CNN	0.0119
A-W2V	Bio	w2v-trained	CNN	0.0145
A-TAGS	Tags	VSM	FF	0.0314
TAGS-ITEMKNN	Tags	-	itemKnn	0.0161
TEXT-RF	Bio	VSM	RF	0.0089
RANDOM	-	-	-	0.0014
UPPER-BOUND	-	-	-	0.5528

# Artist Recommendation: Deep Learning vs. RF

Aproach	Input	Data model	Arch	MAP
A-TEXT	Bio	VSM	FF	0.0161
A-SEM	Sem Bio	VSM	FF	<b>0.0201</b>
A-W2V-GOO	Bio	w2v-pretrain	CNN	0.0119
A-W2V	Bio	w2v-trained	CNN	0.0145
A-TAGS	Tags	VSM	FF	0.0314
TAGS-ITEMKNN	Tags	-	itemKnn	0.0161
TEXT-RF	Bio	VSM	RF	0.0089
RANDOM	-	-	-	0.0014
UPPER-BOUND	-	-	-	0.5528

# Song Recommendation: Audio approach

Approach	Artist Input	Track Input	Arch	MAP
AUDIO	-	audio spec	CNN	0.0015
SEM-VSM	Sem Bio	-	FF	0.0032
SEM-EMB	A-SEM	-	FF	0.0034
<b>MM-LF-LIN</b>	<b>A-SEM</b>	<b>AUDIO emb</b>	<b>MLP</b>	<b>0.0036</b>
MM-LF-H1	<b>A-SEM</b>	AUDIO emb	MLP	0.0035
MM	Sem Bio	audio spec	CNN	0.0014
TAGS-VSM	Tags	-	FF	0.0043
TAGS-EMB	A-TAGS	-	FF	0.0049
RANDOM	rnd emb	-	FF	0.0002
UPPER-BOUND	-	-	-	0.1649

# Song Recommendation: Embeddings vs. VSM

Approach	Artist Input	Track Input	Arch	MAP
AUDIO	-	audio spec	CNN	0.0015
SEM-VSM	Sem Bio	-	FF	0.0032
SEM-EMB	A-SEM	-	FF	0.0034
MM-LF-LIN	<b>A-SEM</b>	<b>AUDIO emb</b>	<b>MLP</b>	<b>0.0036</b>
MM-LF-H1	<b>A-SEM</b>	AUDIO emb	MLP	0.0035
MM	Sem Bio	audio spec	CNN	0.0014
TAGS-VSM	Tags	-	FF	0.0043
TAGS-EMB	A-TAGS	-	FF	0.0049
RANDOM	rnd emb	-	FF	0.0002
UPPER-BOUND	-	-	-	0.1649

# Song Recommendation: Late fusion vs. simultaneous

Approach	Artist Input	Track Input	Arch	MAP
AUDIO	-	audio spec	CNN	0.0015
SEM-VSM	Sem Bio	-	FF	0.0032
SEM-EMB	A-SEM	-	FF	0.0034
<b>MM-LF-LIN</b>	<b>A-SEM</b>	<b>AUDIO emb</b>	<b>MLP</b>	<b>0.0036</b>
MM-LF-H1	<b>A-SEM</b>	AUDIO emb	MLP	0.0035
MM	Sem Bio	audio spec	CNN	0.0014
TAGS-VSM	Tags	-	FF	0.0043
TAGS-EMB	A-TAGS	-	FF	0.0049
RANDOM	rnd emb	-	FF	0.0002
UPPER-BOUND	-	-	-	0.1649

# Song Recommendation: Multimodal vs. single mod.

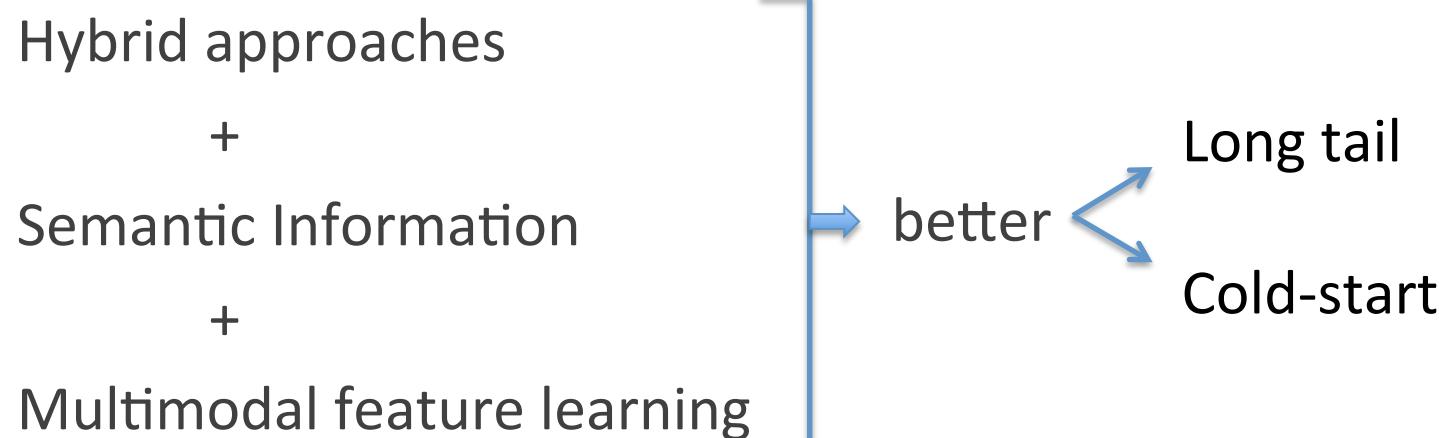
Approach	Artist Input	Track Input	Arch	MAP
AUDIO	-	audio spec	CNN	0.0015
SEM-VSM	Sem Bio	-	FF	0.0032
SEM-EMB	A-SEM	-	FF	0.0034
<b>MM-LF-LIN</b>	<b>A-SEM</b>	<b>AUDIO emb</b>	<b>MLP</b>	<b>0.0036</b>
MM-LF-H1	A-SEM	AUDIO emb	MLP	0.0035
MM	Sem Bio	audio spec	CNN	0.0014
TAGS-VSM	Tags	-	FF	0.0043
TAGS-EMB	A-TAGS	-	FF	0.0049
RANDOM	rnd emb	-	FF	0.0002
UPPER-BOUND	-	-	-	0.1649

# Results

- Splitting the problem between artists and songs ✓
- Learning artist feature embedding separately ✓
- Use of artist biographies ✓
- Semantic enrichment via Entity Linking ✓
- Late fusion of multimodal feature embeddings ✓
- Multimodal better than single modality ✓

# Conclusions

- Knowledge from semantic repositories incorporated via entity linking improves item profiles -> better diversity
- Learning deep feature representations from multimodal data and combining them improves cold-start recommendations



# Reproducibility

- Datasets

MSD-A: <http://www.upf.edu/web/mtg/msd-a>

KG-REC: <https://www.upf.edu/web/mtg/kgrec>

- Source code

TARTARUS: <https://github.com/sergiooramas/tartarus>

lodreclib: <https://github.com/sisinflab/lodreclib>



# Thanks!

@sergiooramas



Universitat  
Pompeu Fabra  
*Barcelona*

**MTG**

Music Technology  
Group