



1. Information in a domain is often spread across multiple (heterogeneous) databases

movies		
id	title	year
m1	The Hangover (2009)	2009
m2	Star Wars: Episode VII - The Force Awakens (2015)	2015

movies2releasedate		
id	month	year
m1	June	2009
m2	December	2015

IMDb

$movies[title] \approx movies2distributors[title] \rightarrow movies[title] \Leftarrow movies2distributors[title]$

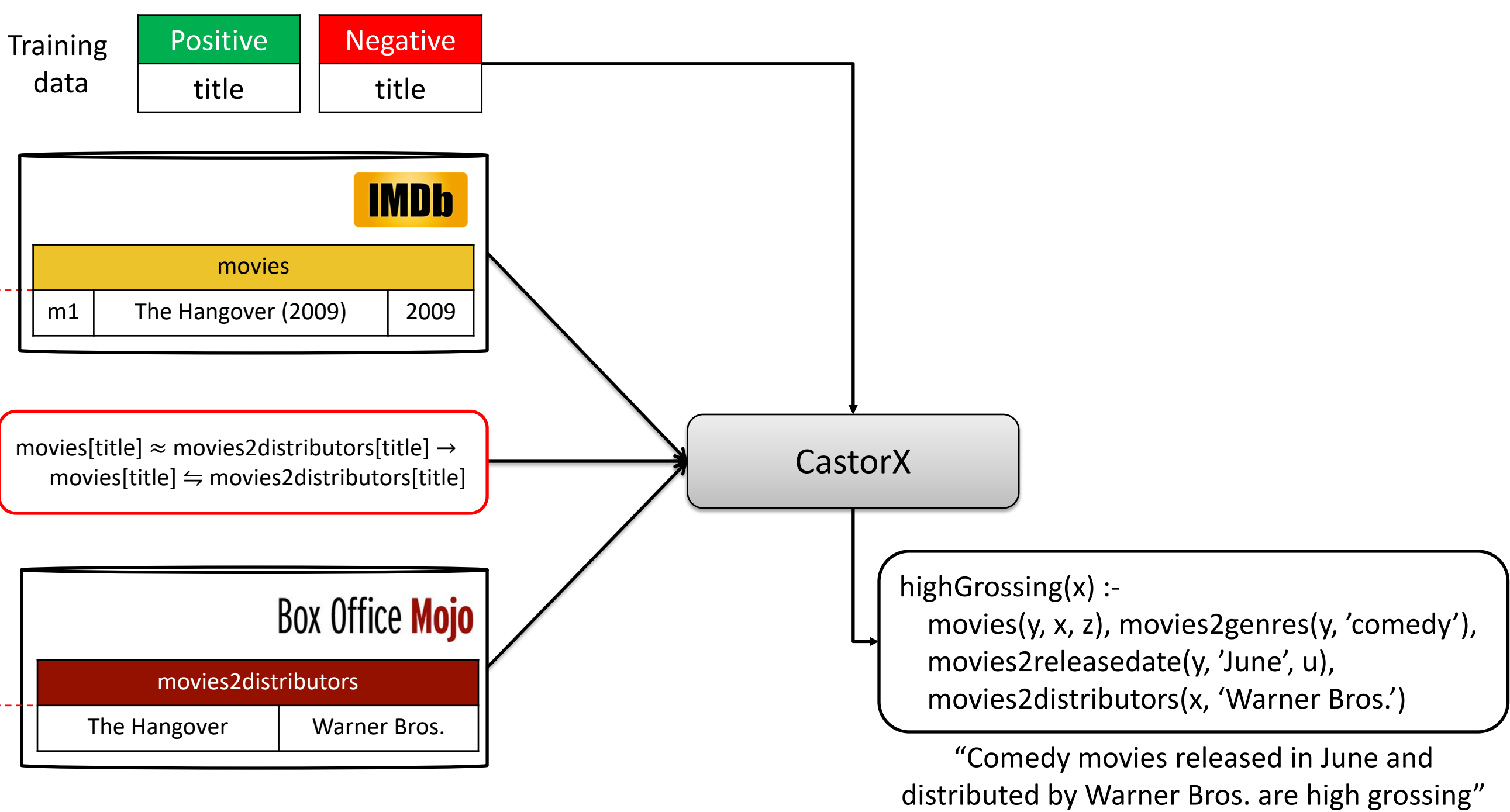
Matching dependencies (MDs) specify relationships between databases.

movies2distributors	
title	distributors
The Hangover	Warner Bros.
Star Wars: The Force Awakens	Buena Vista

movies2budget	
title	budget
The Hangover	35M
Star Wars: The Force Awakens	245M

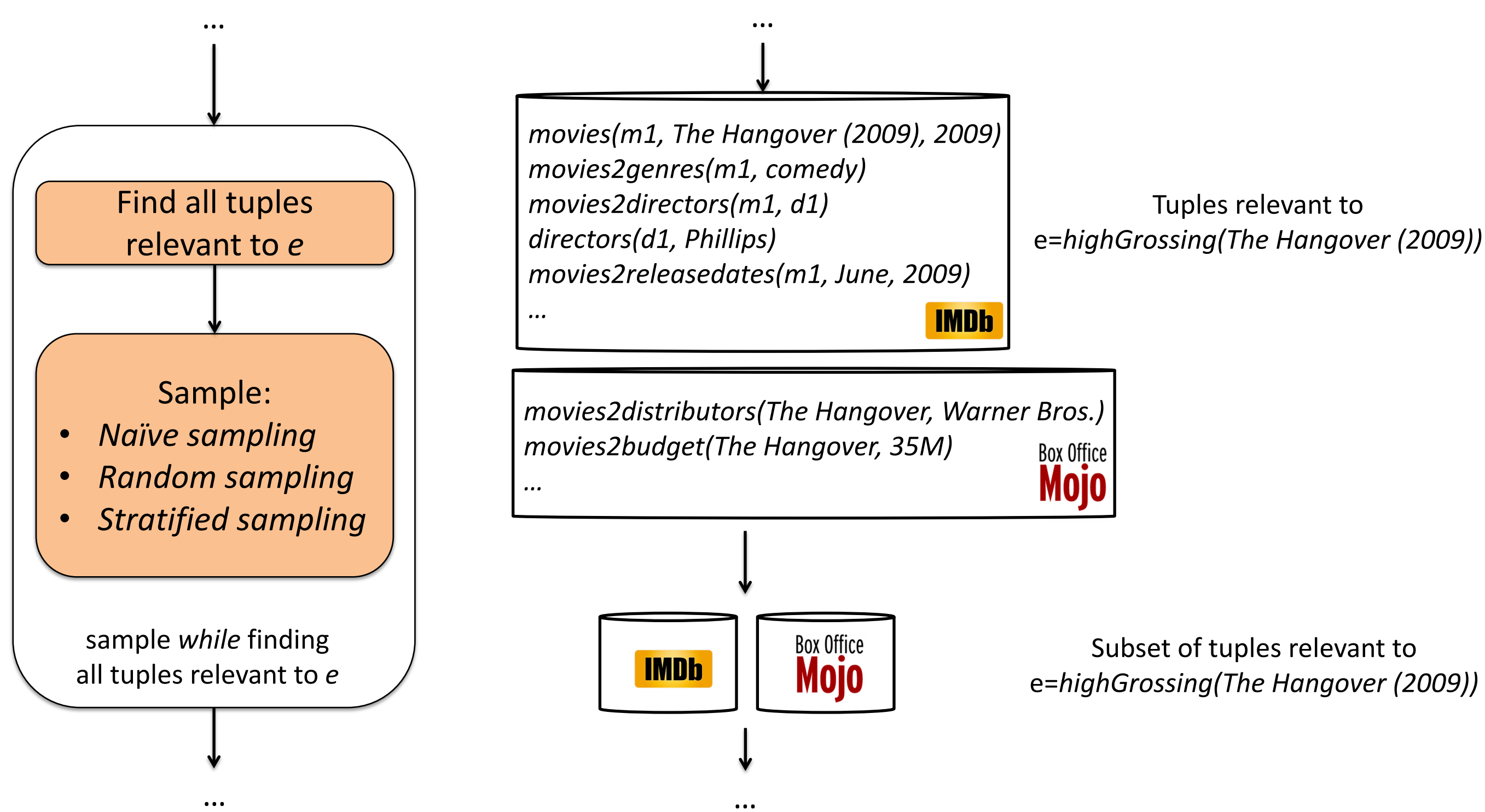
Box Office Mojo

3. CastorX: cross-database relational learning system



5. Efficient learning using sampling

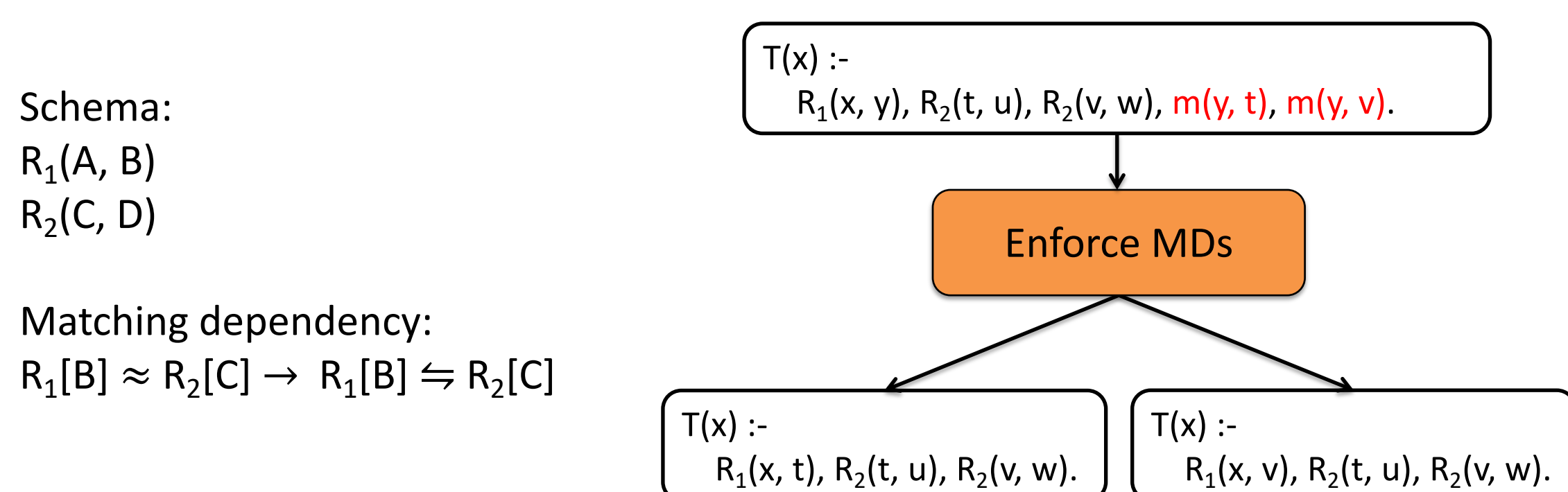
CastorX learns efficiently by sampling the tuples relevant to a positive example e .



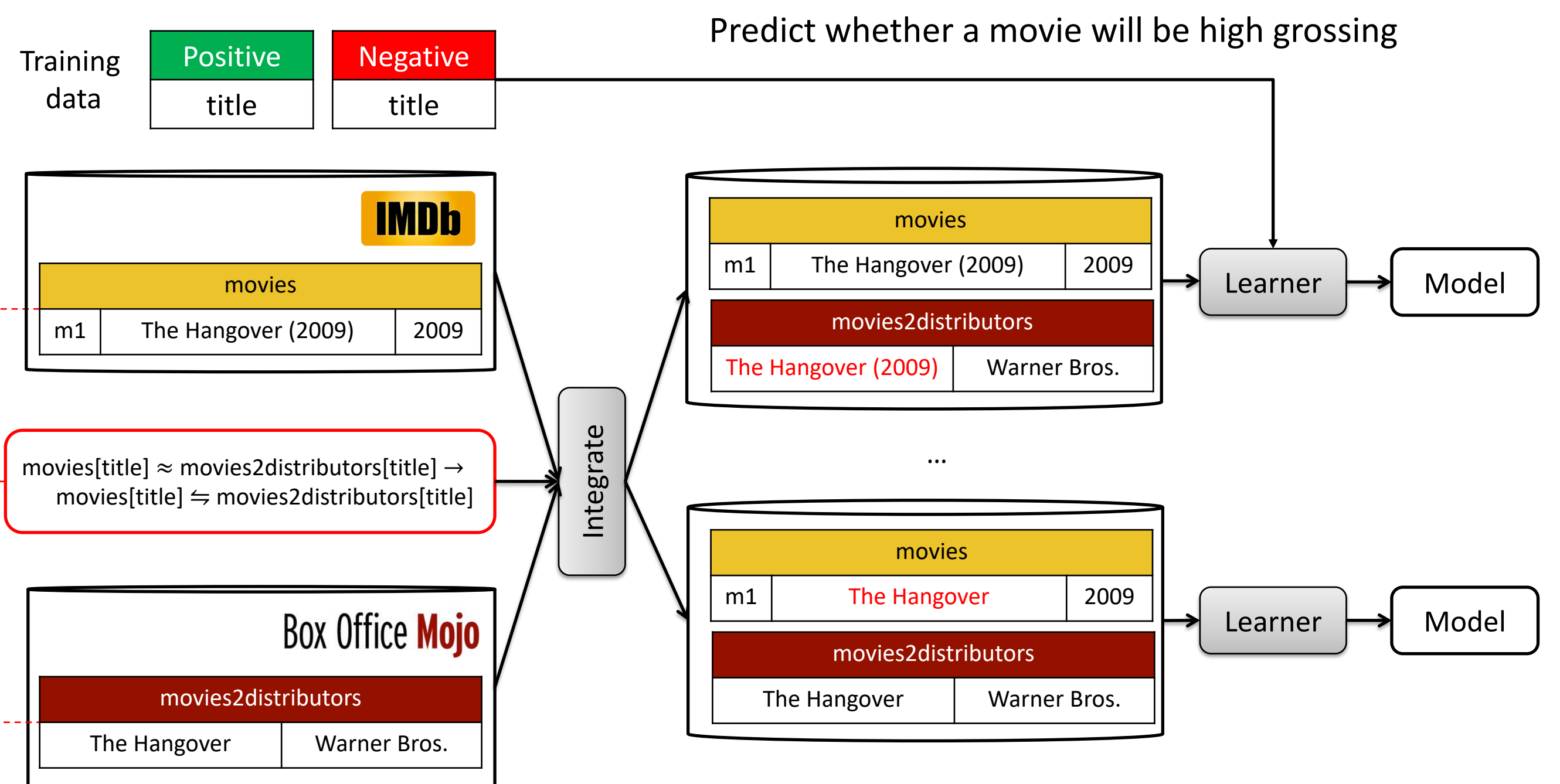
6. Enforcing matching dependencies

After enforcing MDs, learned definitions may be:

- Exactly the same
- Equivalent (through homomorphism)
- Different (can still get approximate answers)



2. Learning over multiple databases

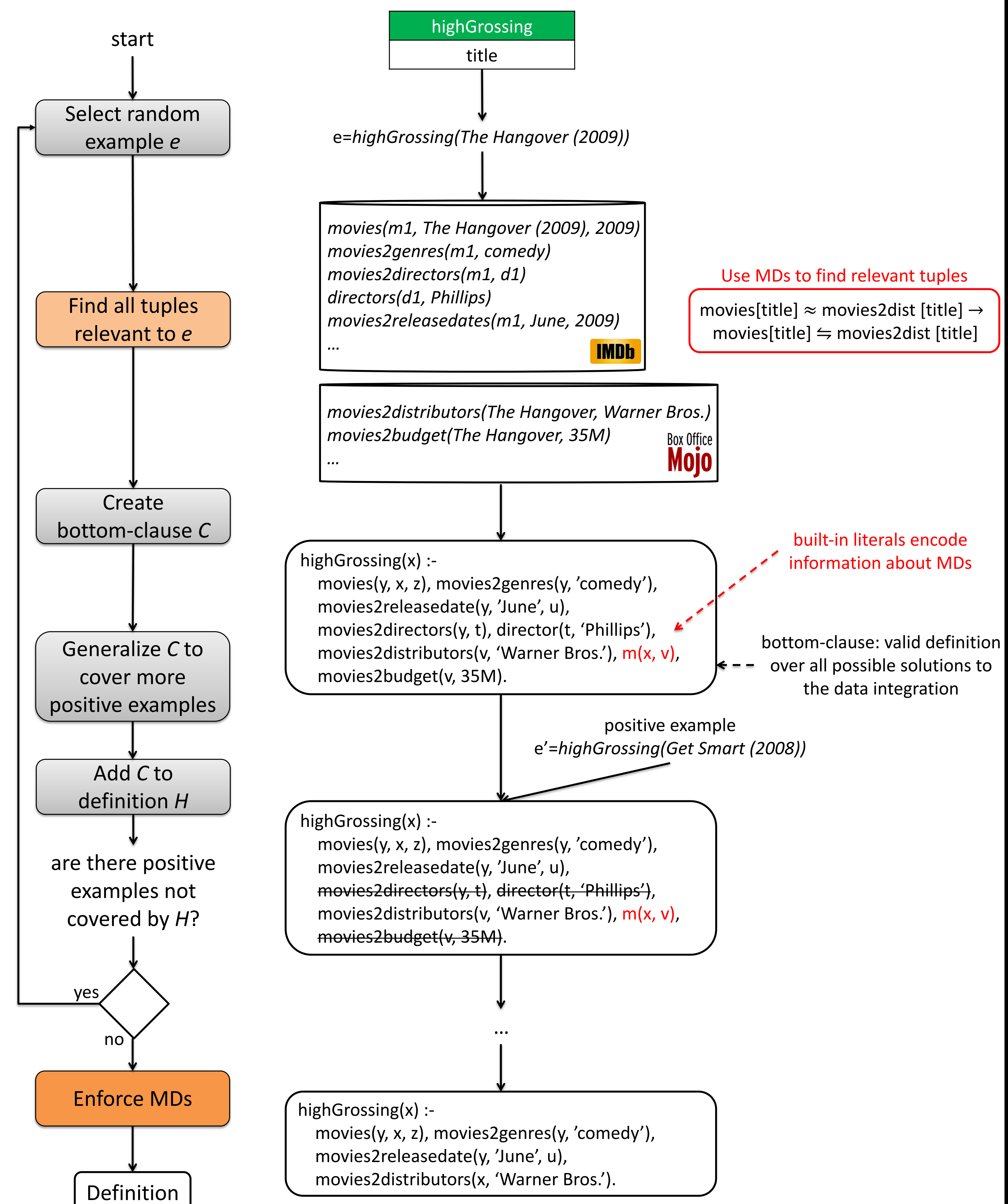


Disadvantages:

- Expensive to integrate databases.
- Creates lots of solutions to the data integration.
- In order to learn, may not need to integrate all data.

4. One representation to encode all possible solutions to the data integration

Encode information about matching dependencies inside the model. Save computations by learning one model that is valid over all possible solutions of the data integration.



7. Experiments

Sampling method	Precision	Recall	Time (minutes)
Naive	0.84	0.87	27.99
Random	0.79	0.81	12.57
Stratified	0.84	0.90	24.97
Naive	0.86	0.78	59.9
Stratified	0.95	0.78	95

HIV DB: chemical compounds

- Target: anti-HIV(compound)
- 7.8M tuples
- 2K positive, 4K negative examples

IMDb + Box Office Mojo

- Target: highGrossing(title)
- 9M and 100K tuples
- 1K positive, 2K negative examples