

# Match-Based Candidate Network Generation for Keyword Queries over Relational Database

Pericles de Oliveira<sup>1,2</sup> Altigran da Silva<sup>1</sup> Edleno de Moura<sup>1</sup> Rosiane Rodrigues<sup>1</sup>

NOKIA Solutions and Networks, Rio de Janeiro, Brazil<sup>2</sup>



## Introduction

- ▶ R-KwS systems: take queries as a set of keywords and return JNTs – joint networks of tuples – that fulfill the user needs
- ▶ Candidate Networks (CNs): Relational Join Expressions automatically generated from input keywords that when evaluated by a RDBMS produce relevant JNTs
- ▶ Current CN generation approaches require an exhaustive exploration of a large number of combinations of keyword occurrences in the DB
- ▶ For certain queries, they can take too long to produce answers and may even fail to return results (e.g., by exhausting memory)
- ▶ We propose MatCNGen, a novel approach for generating CNs that drastically reduces the time required to generate CNs
- ▶ Besides improving CN generation, MatCNGen also has a positive impact on the evaluation of CNs: run faster and produces better results.

## MatCNGen

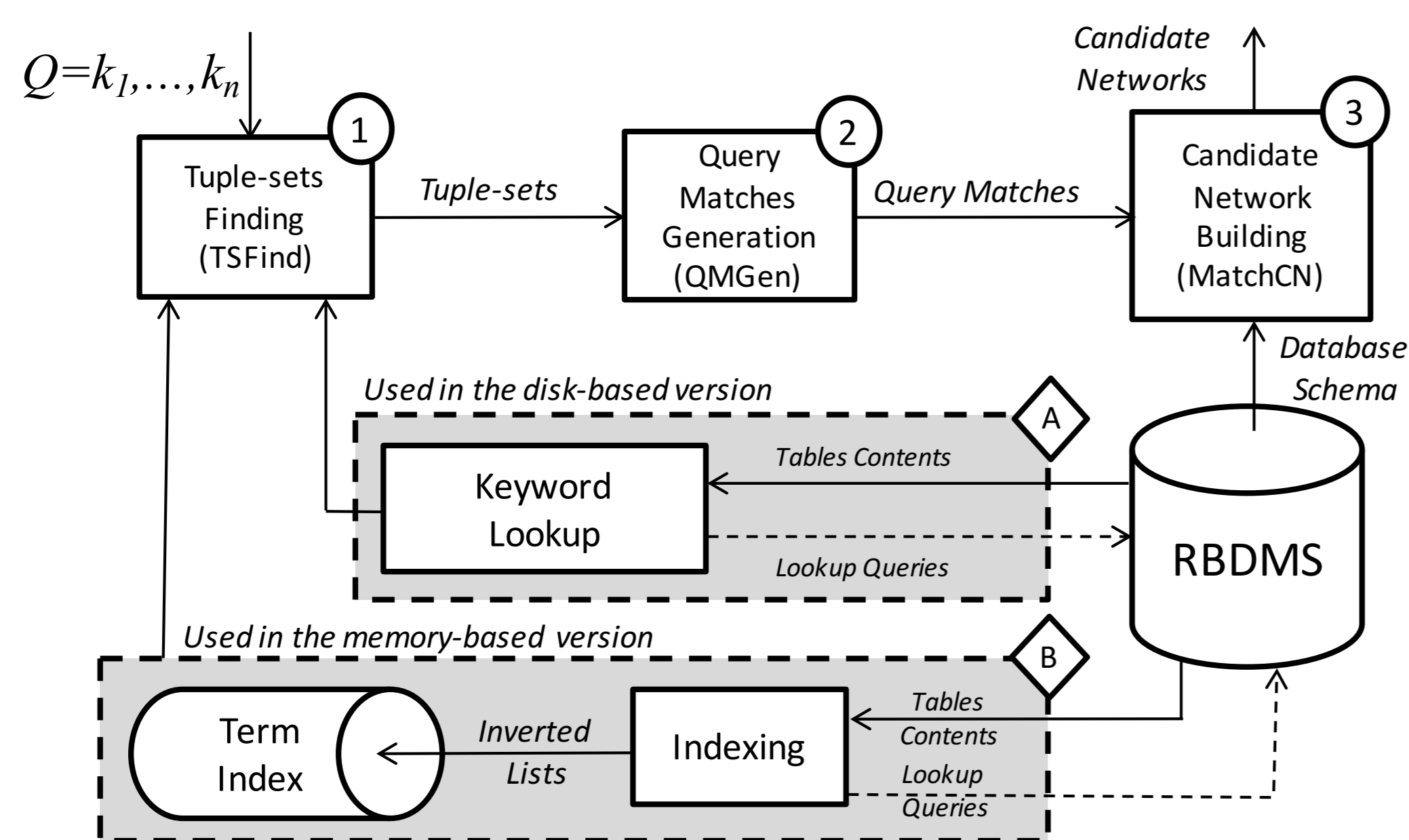
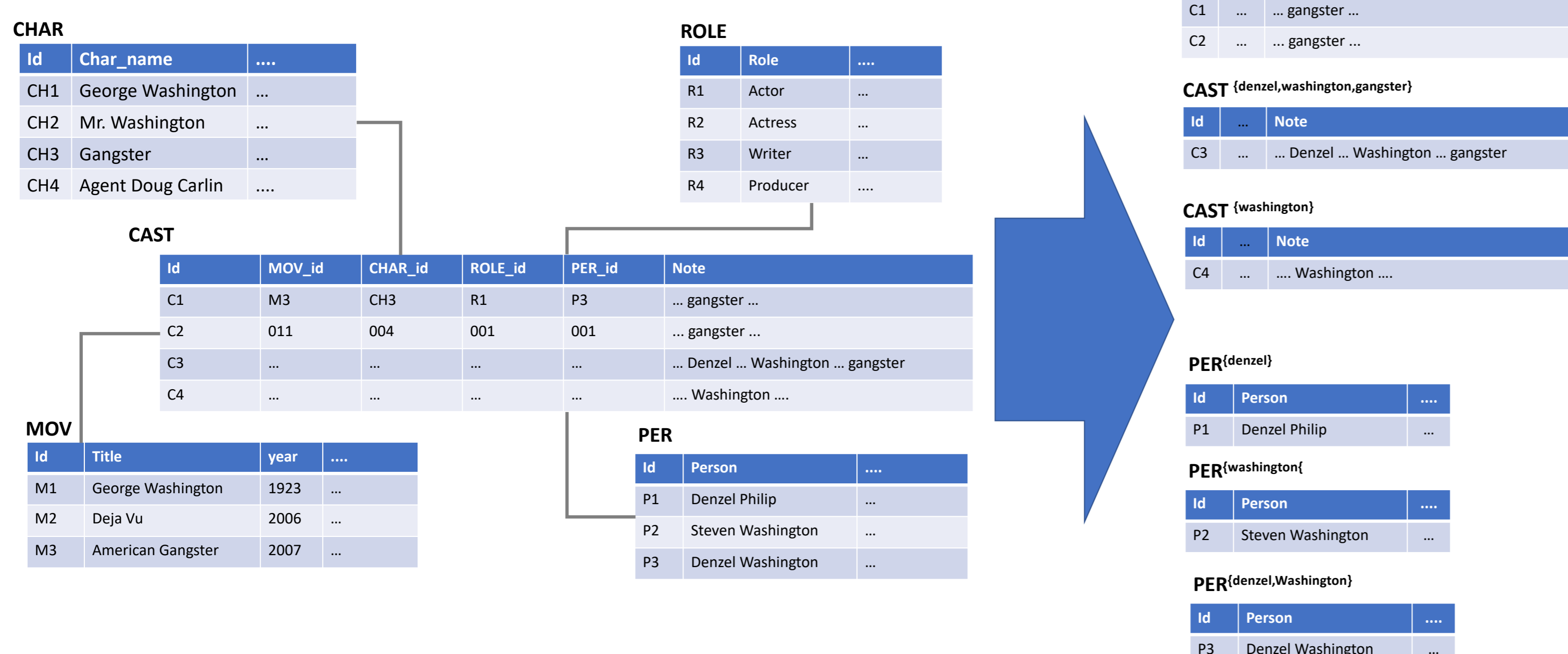


Figure: MatCNGen main steps.

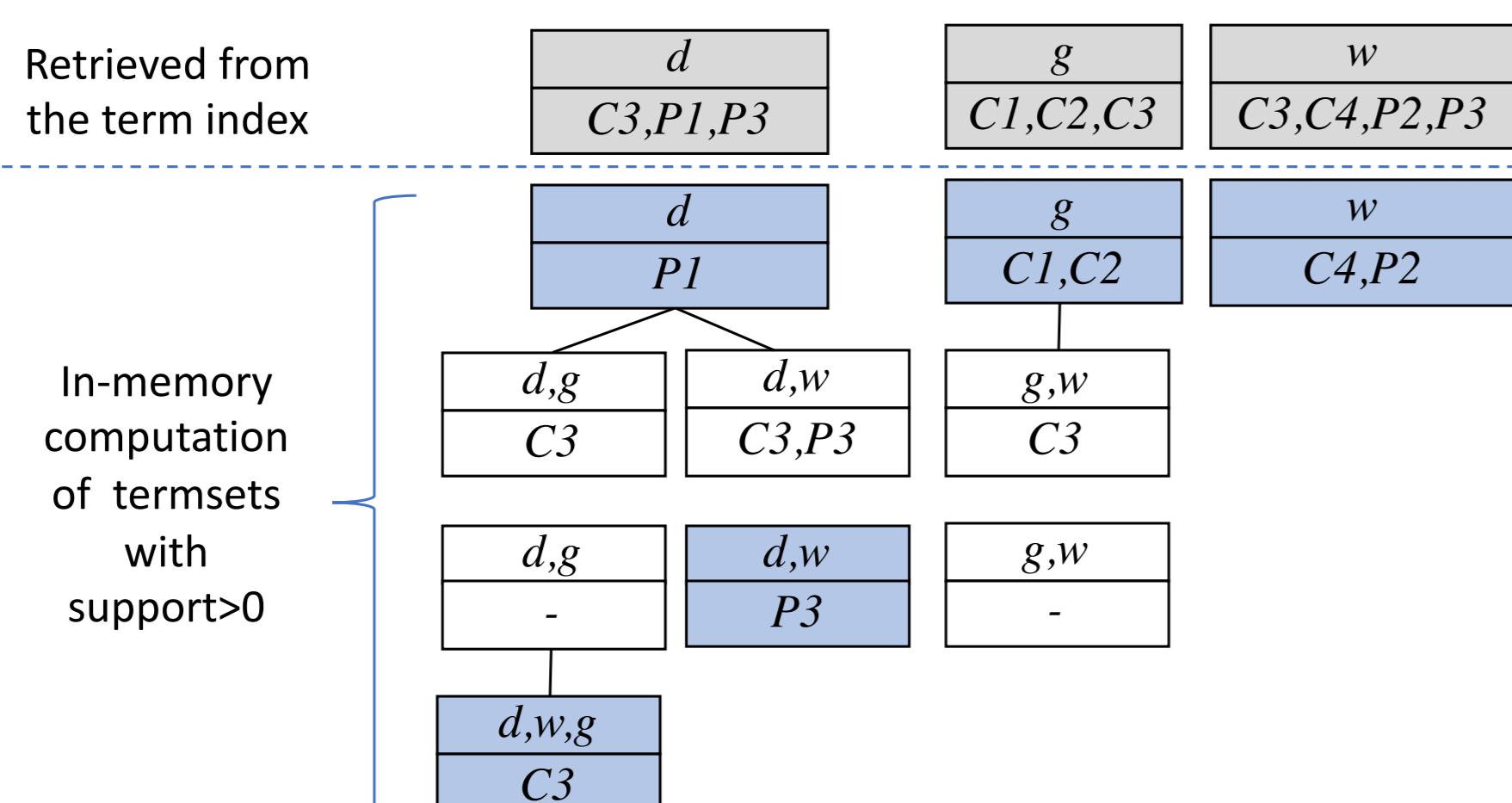
## Step 1 – Finding Tuple Sets

What are the relevant pieces of information in the DB?



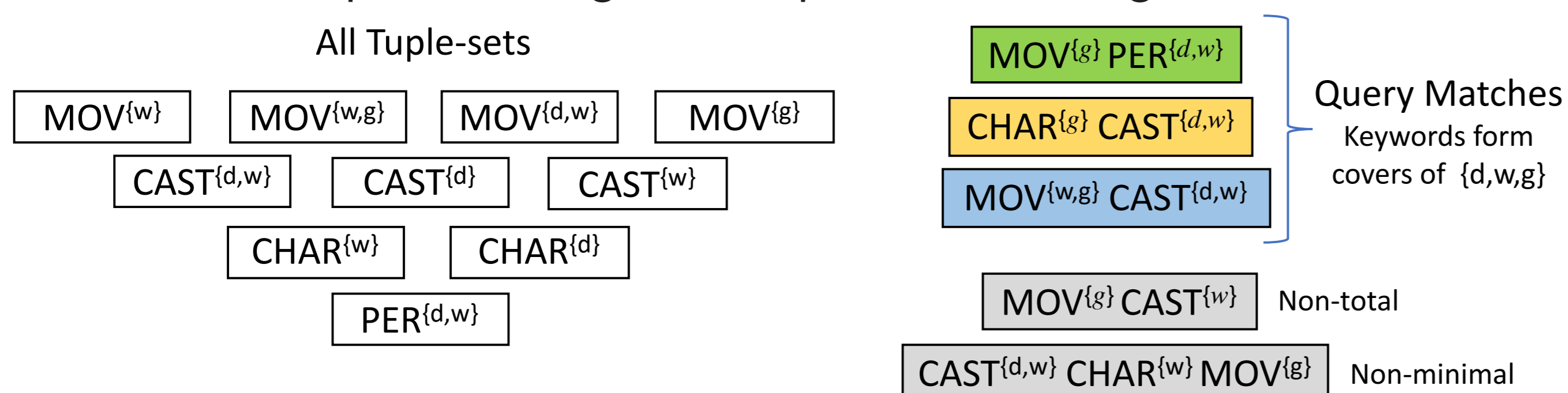
Database Instance

Tuple-Sets



## Step 2 – Query Matches Generation

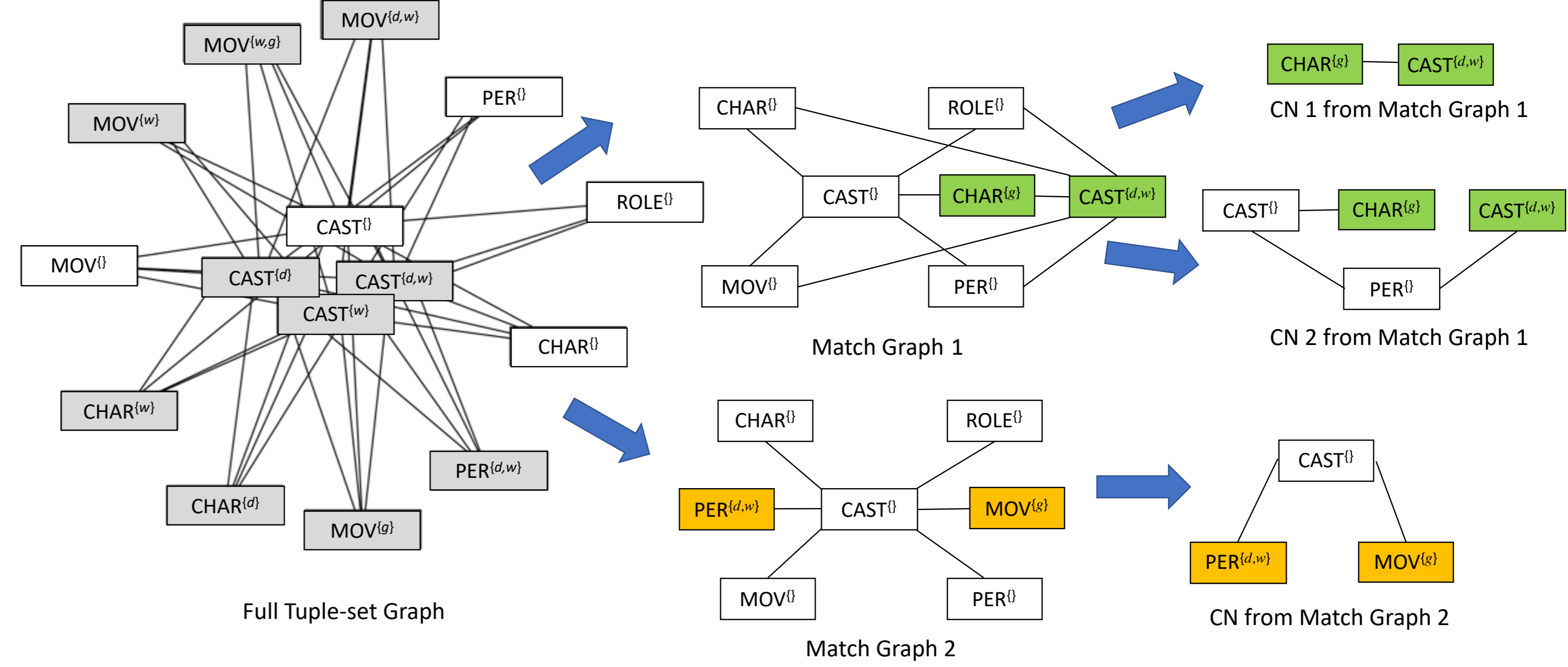
Which pieces fit together to produce meaningful answers?



- ▶ For a query  $Q$ , there is no Query Match with more than  $|Q|$  tuple-sets!!

## Step 3 – Candidate Network Building

How to join the pieces that fit together?



$\sigma_{\text{title} \supseteq \{\text{gangster}\}} \text{MOV} \bowtie_{\text{id}=\text{cid}} \text{CAST} \bowtie_{\text{cid}=\text{pid}} \sigma_{\text{name} \supseteq \{\text{denzel, washington}\}} \text{PER}$

## Experimental Datasets

Dataset	Size (MB)	Relations	Tuples	RIC
Mondial	9	28	17,115	104
IMDb	516	5	1,673,074	4
Wikipedia	550	6	206,318	5
DBLP	40	6	878,065	6
TPC-H	876	8	2,389,071	11

Table: Characteristics of the datasets used.

## Results

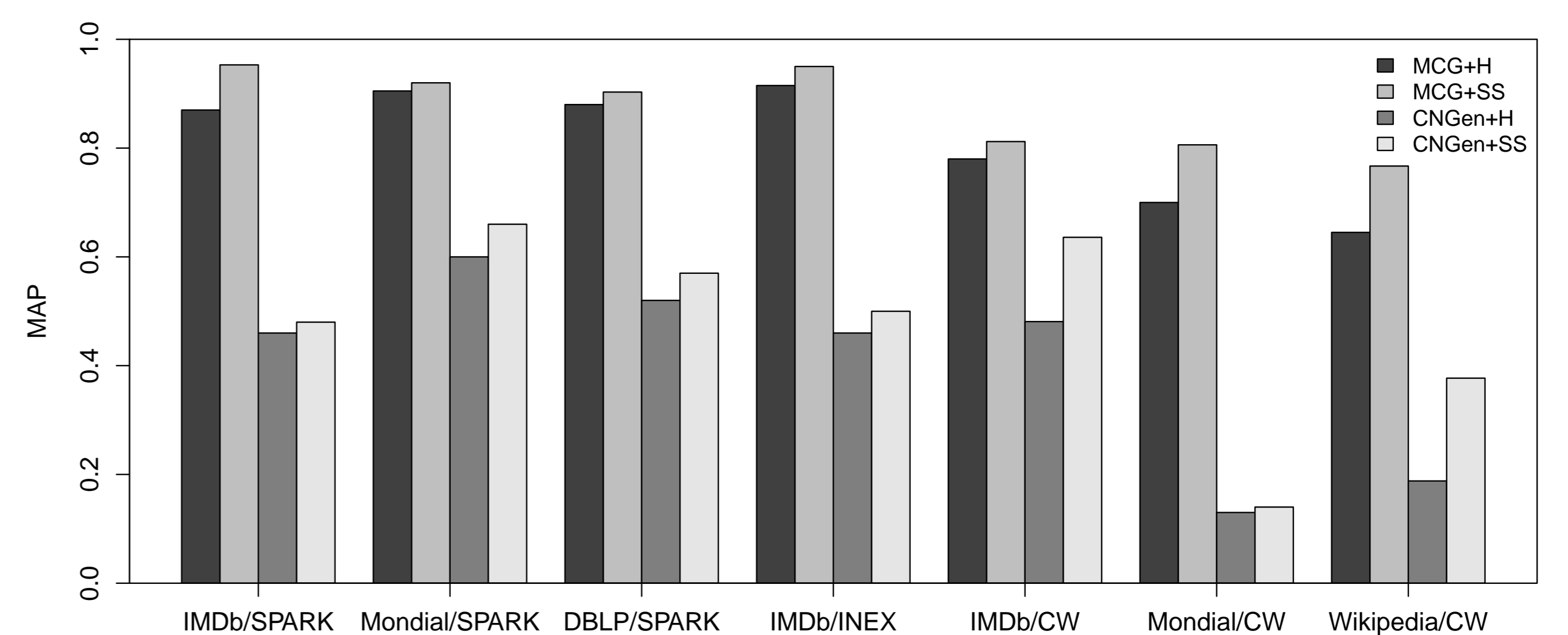


Figure: MAP results for all query sets.

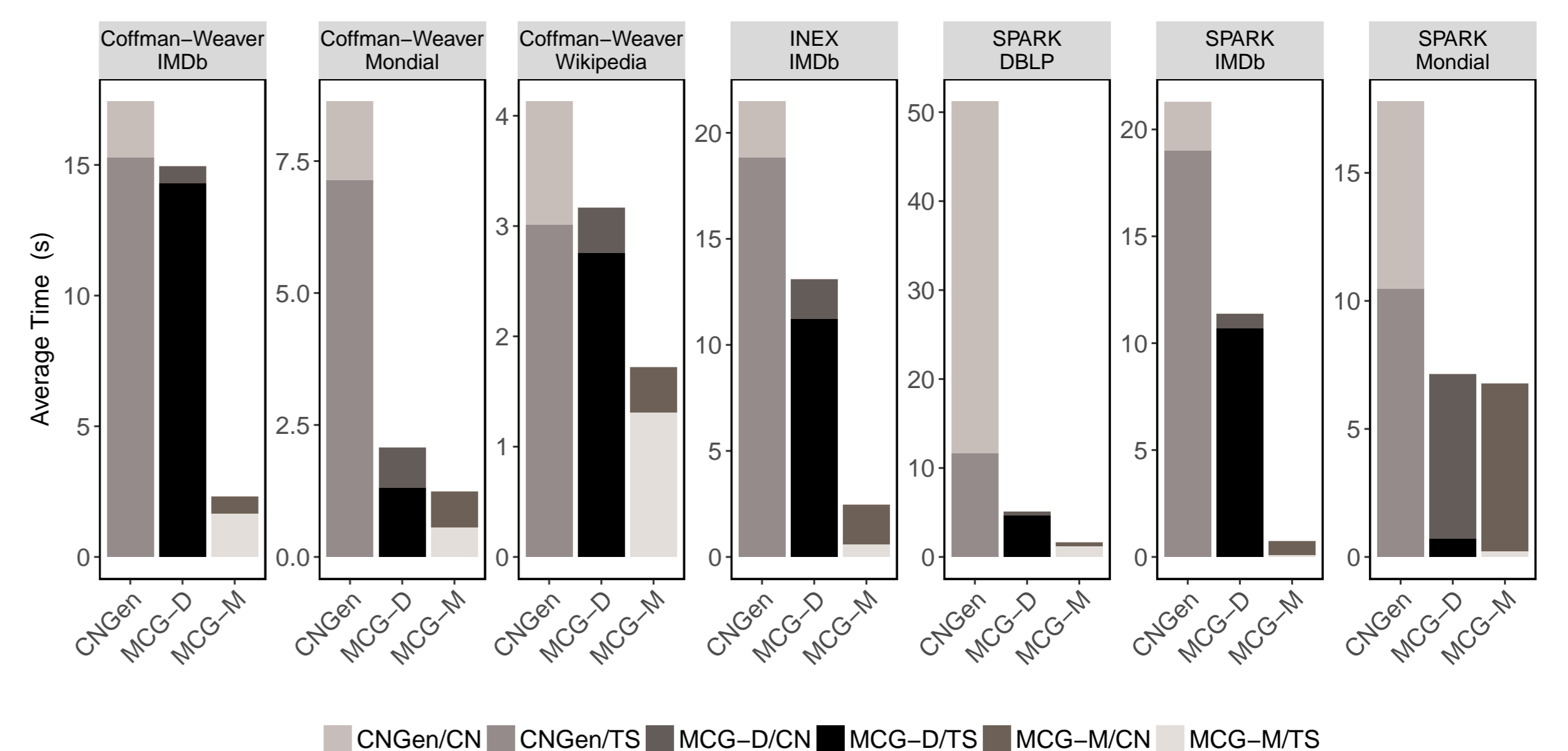
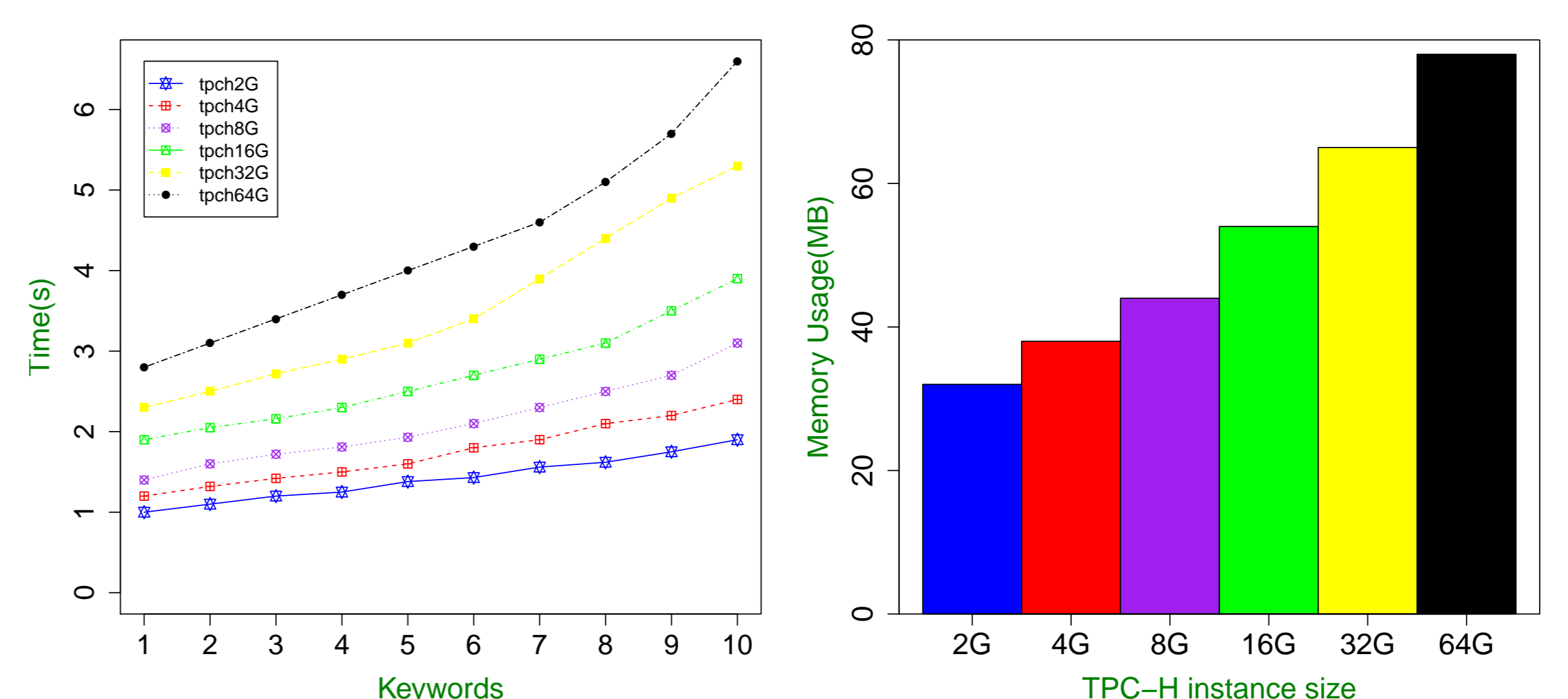


Figure: Time to generate CNs using CNGen and MatCNGen.



Response Time (a)

Memory Usage (b)

Figure: Scalability with instances of varying sizes.