

Learning to Query: Focused Web Page Harvesting for Entity Aspects

Yuan Fang¹ Vincent Zheng² Kevin Chang^{2,3}

¹ Institute for Infocomm Research, Singapore

² Advanced Digital Sciences Center, Singapore

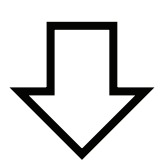
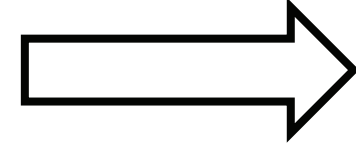
³ University of Illinois at Urbana-Champaign, USA



Problem: Learning to Query (L2Q)

Entity type	Common aspects
celebrity	spouse, age, net worth
car	safety, cost, interior
business	address, hours, phone

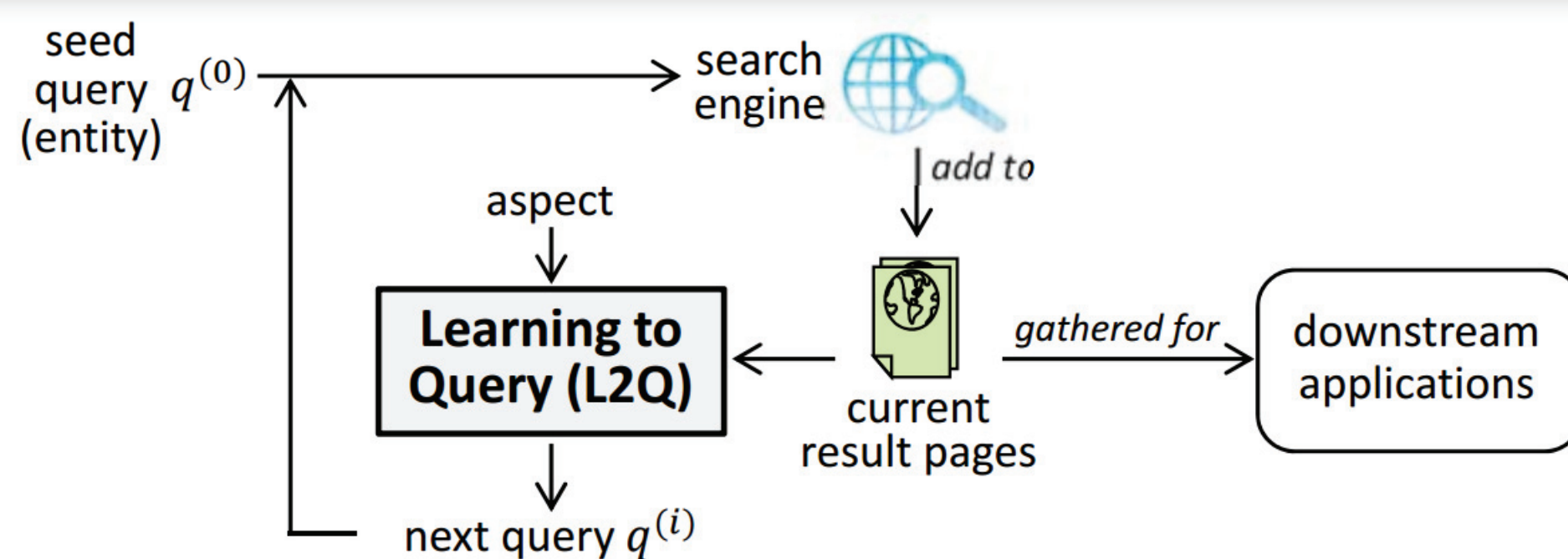
abundant but
scattered



Business
Analytics

Vertical
portal or
search

Overall Workflow: Iterative Querying



Seed query:

Keywords (uniquely)
identifying the entity

Target aspect:

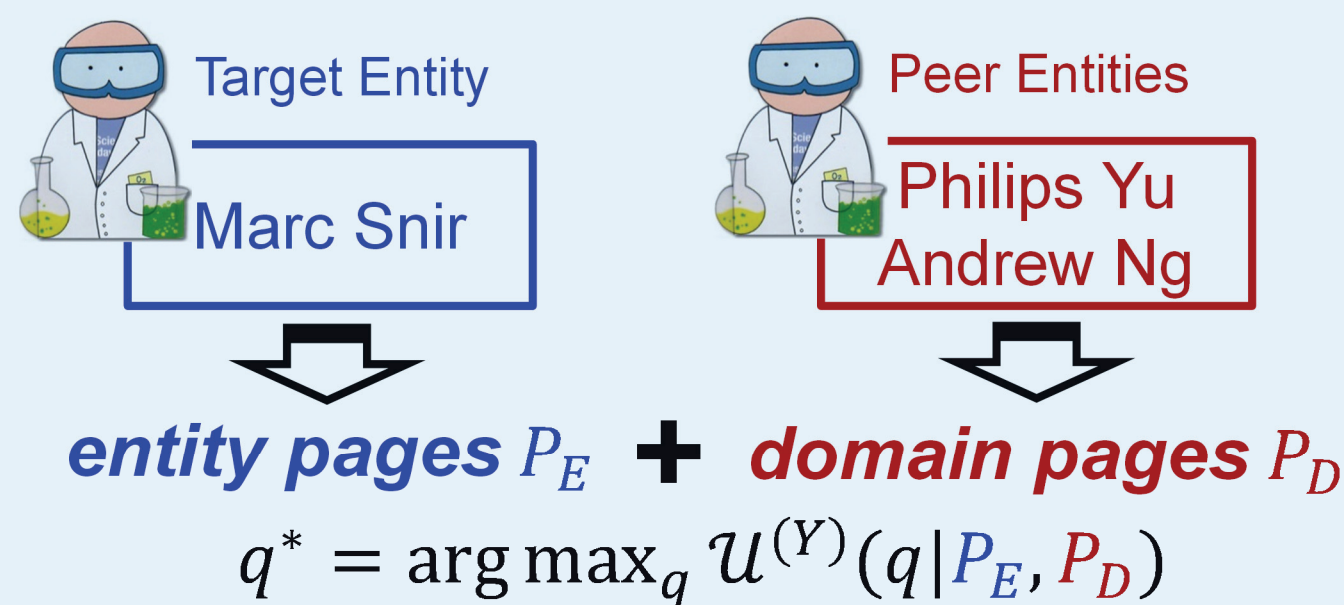
A pre-trained classifier Y
for the target aspect

**Utility:
(precision/recall)**

In each iteration,
 $q^* = \arg \max_q U^{(Y)}(q)$

Subproblem #1: Domain-aware L2Q

a) Domain Pages

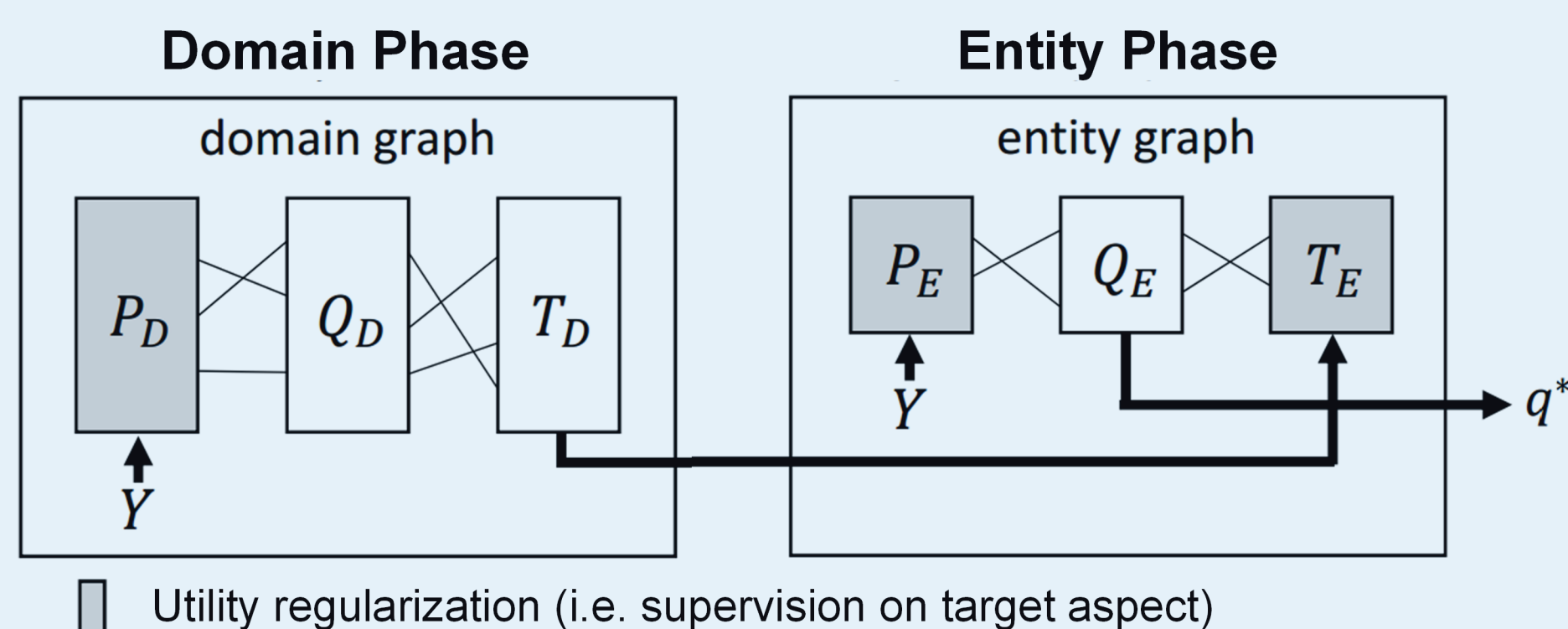


b) Template Abstraction

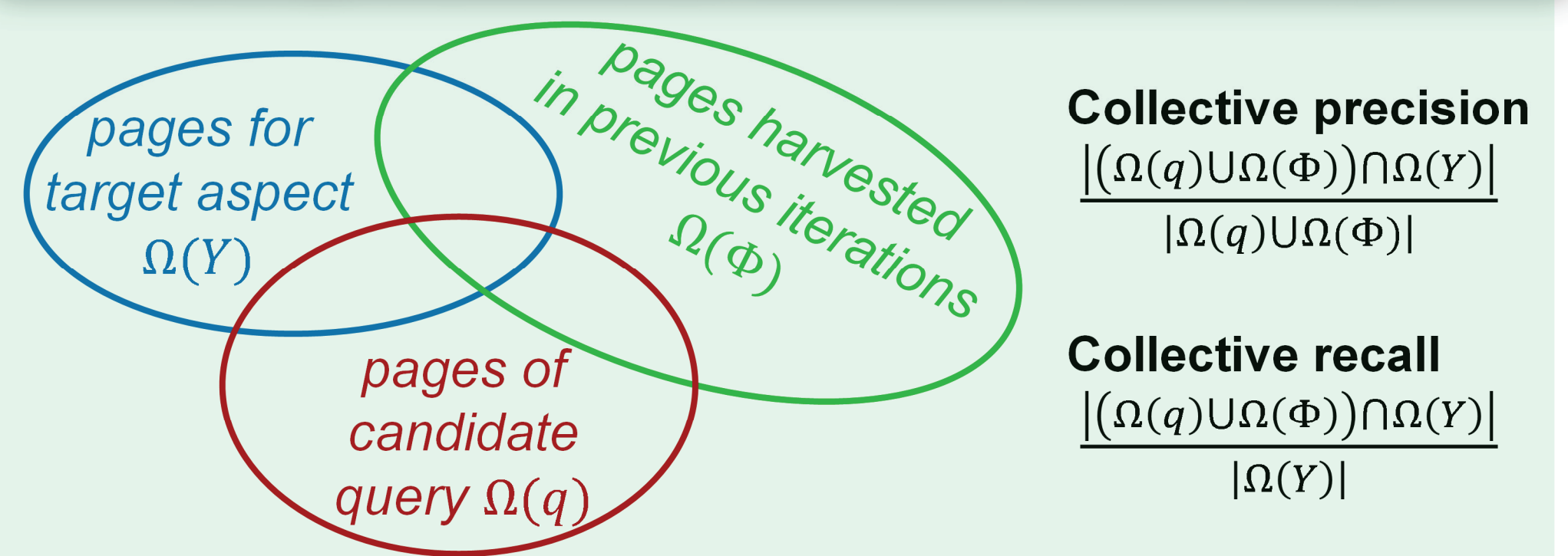
Entity	Example page content	Example query
Marc Snir	...many HPC papers in IJHPCA ...	hpc ijhpca
Philip Yu	...his data mining papers in TKDE ...	data mining tkde
Andrew Ng	...his recent AI paper in JMLR ...	ai jmlr

<topic> <journal>

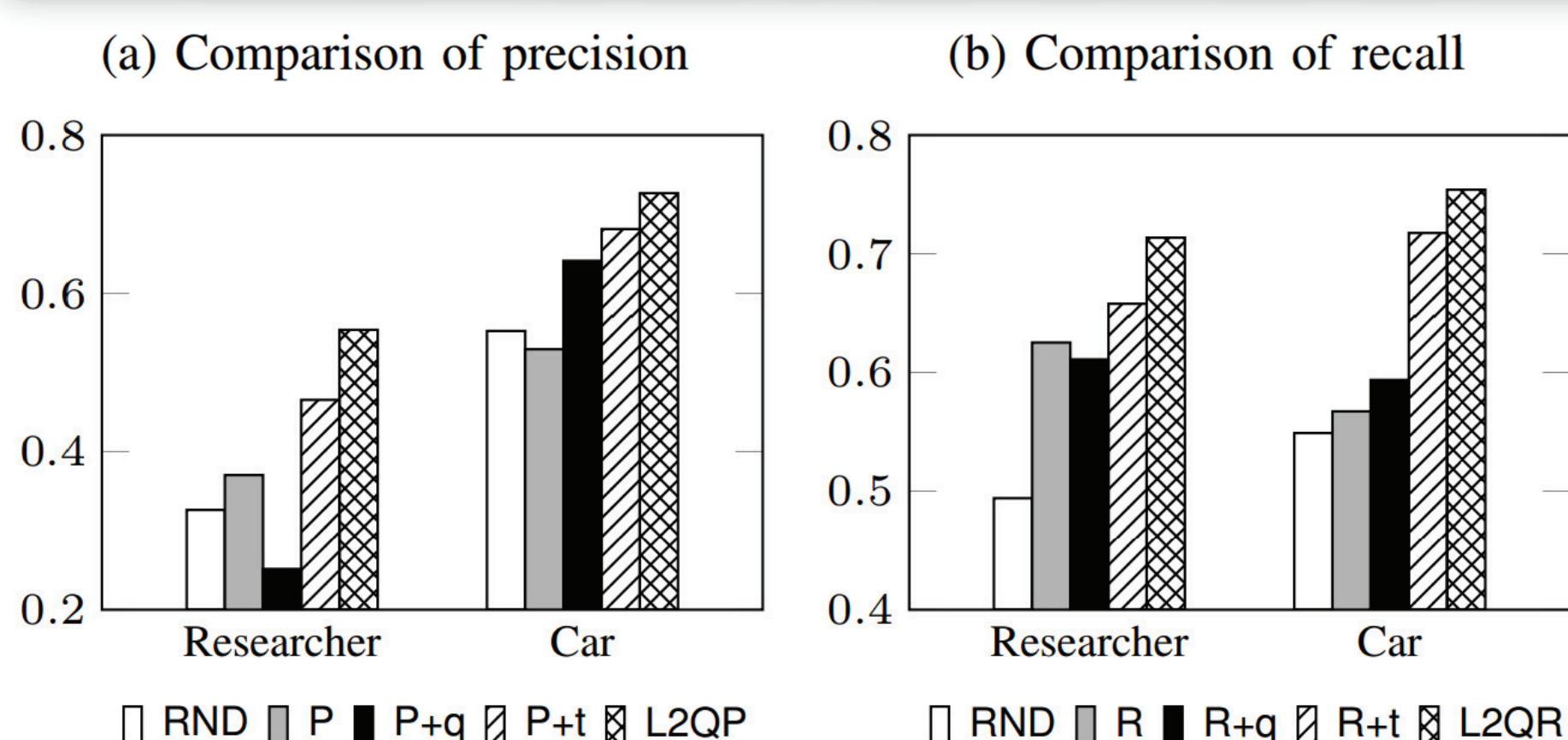
c) Bridging Domain & Entity



Subproblem #2: Context-aware L2Q



Result A: Effect of Domain+Context



RND: select query randomly

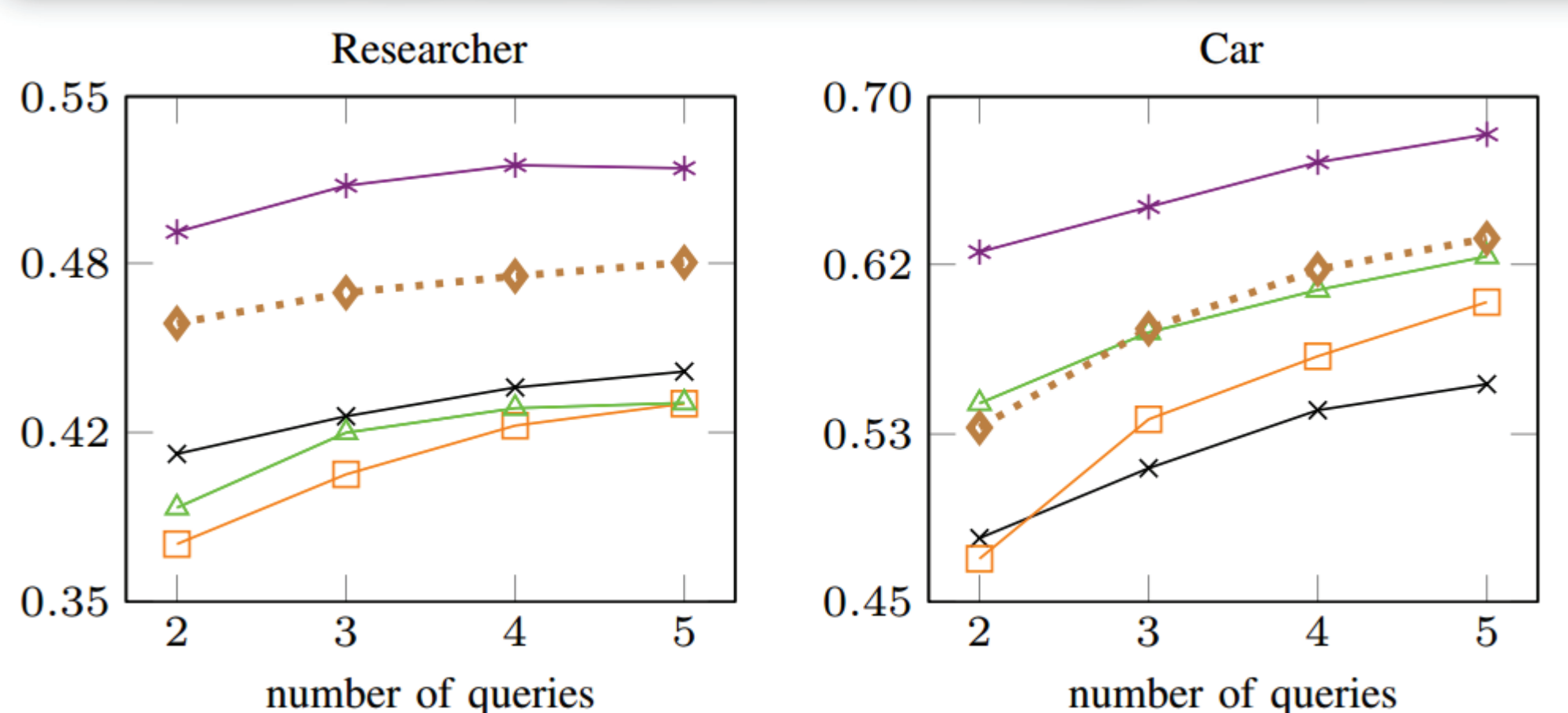
P/R: optimize precision/recall without domain and context

P/R+q: with domain pages but no templates, and without context

P/R+t: with domain pages and templates, without context

L2QP/L2QR: full approaches optimizing precision/recall

Result B: Compare with Indep. Baselines



L2QBAL: optimize for F-score, balancing L2QP & L2QR

LM: language feedback model

AQ: adaptive querying for text databases

HR: harvest rate for hidden structured databases

MQ: manually designed queries

—*— L2QBAL

—x— LM

—□— AQ

—△— HR

—◇— MQ