



PQPP: A Joint Benchmark for Text-to-Image Prompt and Query Performance Prediction

Eduard Poesina¹, Adriana-Valentina Costache¹, Adrian-Gabriel Chifu², Josiane Mothe³, Radu Tudor Ionescu¹

¹University of Bucharest ² Aix-Marseille Université, LIS ³ Univ. de Toulouse, Institut de Recherche en Informatique de Toulouse



Introduction

Prompt Performance Prediction – Pre/Post-Generative



Training Data



Testing Data









Introduction

Github



Query Performance Prediction – Pre/Post-Retrieval



Training Data







Introduction

Context & Challenges:

No prior work explores the relationship between Query Performance Prediction (QPP) and Prompt Performance Prediction (PPP)

Research questions:

RQ1: Are the two tasks related, or they require specifically trained models? **RQ2**: Could a QPP model also be used for PPP and vice-versa?

Our contributions:

1. The first joint benchmark for query/prompt performance prediction

- 2. A database of ~ 10K queries/prompts, with > 1.5M total annotations
- 3. Baselines: multiple pre-generative/retrieval and post-generative/retrieval





Related Work

Prompt Performance Prediction

Bizozerro et al,2024 - Prompt Performance Prediction and training models on automated relevance scoring
Pavlichenko et al., 2023, SIGIR – Method to adapt a given prompt to improve text to image model performance
Kirstain et al., 2023, NeurIPS - Preference annotations by asking humans to choose an image from a pair of generated images

Query Performance Prediction

Xing et al,2010,ECIR - Groundwork of exploring query difficulty prediction in image retrieval

Liu et al., 2012, Neurocomputing – Focus on estimating query difficulty with unigram language models and visual word verification





Proposed Benchmark

- Source datasets for prompts:
 - MS COCO: 10K captions (sampled via k-means for diversity)
 - DrawBench: all 200 prompts (more suitable for generation)
- Generative models:
 - GLIDE: two image samples per prompt
 - SDXL: two image samples per prompt
- Retrieval Models:
 - CLIP
 - BLIP2





Prompt Performance Assessment

Men sit in the outfield in white uniforms waiting for the pitch.



high relevance
 low relevance
 no relevance
 unrealistic



high relevance
low relevance
no relevance
unrealistic



high relevance
 low relevance
 no relevance
 unrealistic



high relevance
low relevance
no relevance
unrealistic



high relevance
low relevance
no relevance
unrealistic

Label definitions:

- high relevance over half of the concepts mentioned in the prompt
- low relevance at least one concept, but fewer than half
- no relevance unrelated, yet realistic
- Unrealistic notable artifacts





Prompt Performance Assessment

Target Annotations: 10,200 prompts × 5 images × 3 annotations = 153,000 annotations

Criteria for exclusion:

Cohen's K agreement with control prompts < 0.4

Annotators:

Recruited: 173 Validated: **147 (247,050 annotations)**

Statistic	Min	Mean	Max
#annotations per person	30	1,681	15,845
Fleiss' κ	0.41	0.54	1.00

Table 1. Statistics about the annotators enrolled in the annotation process for generated images.

https://www.researchgate.net/figure/Fleiss-Kappa-and-Inter-rater-agreement-interpretation-24_tbl3_281652142





Query Performance Assessment





Query Performance Assessment

Target Annotations: 10,200 prompts × Max 2000 images × 3 annotations Annotated Images: 1,393,363

Criteria for exclusion: F1 score < 0.4 with control prompts

Annotators:

Recruited: 100

Validated: 93

Min F_1	Mean F ₁
0.477	0.727



Predictors

Github



Pre-generation/retrieval Predictors

Basic text predictors

diversity of concepts lexical density morphological complexity frequency of grammatical structures

Fine-tuned BERT Base architecture, based on cased inputs

Post-generation/retrieval Predictors

Fine-tuned CLIP Long-CLIP with a Vit-B/32 backbone

Correlation-based CNN Model trained over the correlation matrix between image pairs





Experiments and Results

Predictor Type	Predictor Name	Generative Task			Retrieval Task								
		GLIDE		SDXL		CLIP			BLIP-2				
		HBPP		HBPP		P@10		RR		P@10		RR	
		rson	Idall	rson	idall	rson	Idall	rson	Idall	rson	Idall	rson	Idall
		Pea	Ken	Pea	Ken	Pea	Ken	Pea	Ken	Pea	Ken	Pea	Ken
Post- Pre-	#synsets	-0.112^{\ddagger}	-0.076^{\ddagger}	-0.087^{\ddagger}	-0.080^{\ddagger}	-0.110^{\dagger}	-0.058^{\ddagger}	-0.034	-0.012	-0.115^{\ddagger}	-0.070^{\ddagger}	-0.038	-0.010
	#words	$ -0.090^{\dagger} $	$ -0.084^{\ddagger} $	$ -0.105^{\ddagger} $	-0.109^{\ddagger}	-0.133^{\ddagger}	$ -0.104^{\ddagger} $	-0.035	-0.026	$ -0.175^{\ddagger} $	-0.136^{\ddagger}	-0.038	-0.015
	Fine-tuned BERT	0.566^{\ddagger}	0.406^{\ddagger}	0.281^{\ddagger}	0.232^{\ddagger}	0.451^{\ddagger}	0.277^{\ddagger}	0.221‡	0.176‡	0.511 [‡]	0.328^{\ddagger}	0.168^{\ddagger}	0.139^{\ddagger}
	Fine-tuned CLIP	0.649 [‡]	0.474^{\ddagger}	0.380^{\ddagger}	0.246^{\ddagger}	0.473^{\ddagger}	0.299 [‡]	0.200^{\ddagger}	0.149^{\ddagger}	0.498^{\ddagger}	0.358^{\ddagger}	0.166^{\ddagger}	0.150^{\ddagger}
	Correlation CNN	0.548^{\ddagger}	0.393^{\ddagger}	0.159^{\ddagger}	0.107^{\ddagger}	0.270^{\ddagger}	0.186^{\ddagger}	0.189^{\ddagger}	0.162^{\ddagger}	0.159^{\ddagger}	0.133^{\ddagger}	0.206‡	0.158^{\ddagger}

Predictor Results





Experiments and Results



Metric	Pearson	Kendall
HBPP vs. P@10	0.135 [‡]	0.093 [‡]
HBPP vs. RR	0.072‡	0.048†
P@10 vs. RR	0.560 [‡]	0.512 [‡]







Experiments and Results



LIDE	Predicted Score	SDXL
+	0.67	
	0.49	
1	-0.02	KK
	0.01	
	-0.23	

Predicted

Score

1.71

2.00

1.87

1.95

1.60

Original Captions	Rephrased Captions					
An older man is holding a surfboard while a young boy stands on it. 1.88	A young boy standing on a surfboard, held steady by an elderlyman. 1.95	An older man grips a surfboard, supporting a young boy standing on it. 1.88	An elderly man steadying a surfboard as a young boy balances on top. 1.83			
			the set			
A study table with computer, mouse and keyboard. Photo frame are also kept.	A neatly arranged study table with a computer, mouse, keyboard, and a photo frame nearby.	A study desk featuring a computer setup with a mouse and keyboard, alongside a photo frame.	A computer, mouse, and keyboard on a study table, with a photo frame placed beside them.			
1.70	1.90	1.79	1.74			
Two little giraffes standing between two slightly bigger ones. 0.30	Two small giraffes standing between two slightly taller giraffes. 1.48	Two young giraffes nestled between two slightly bigger ones. 1.08	A pair of little giraffes positioned between two larger giraffes.			



Thank you!

PQPP: A Joint Benchmark for Text-to-Image Prompt and Query Performance Prediction

Eduard Poesina¹, Adriana-Valentina Costache¹, Adrian-Gabriel Chifu², Josiane Mothe³, Radu Tudor Ionescu¹

¹University of Bucharest ² Aix-Marseille Université, LIS ³ Univ. de Toulouse, Institut de Recherche en Informatique de Toulouse



