

# Automatically Extracting Numerical Results from Randomized Controlled Trials with LLMs

Hye Sun Yun, David Pogrebitskiy, Iain J Marshall, Byron C Wallace

{yun.hy, pogrebitskiy.d,b.wallace}@northeastern.edu, {iain.marshall}@kcl.ac.uk



## MOTIVATION

**Meta-analyses** of randomized controlled trials (RCTs) provide robust estimates of treatment efficacy and require extraction of data elements from individual articles for synthesis.

- Can we fully automate “on-demand” meta-analysis of evidence relevant to a given clinical question?
- Are modern LLMs sufficiently capable of *numerical* data extraction to permit accurate, fully automated meta-analysis?

## DATA ANNOTATION

- **Intervention, Comparator, & Outcome (ICO triplets)** from PubMed RCT reports
- Annotations based on **Abstract + Results** sections of RCT
- Schema:
  - **Type of outcome:** binary or continuous
  - **Binary outcome:** events, group sizes for I & C
  - **Continuous outcome:** means, standard deviations, group sizes for I & C

### Example Annotation for Given ICO Triplet

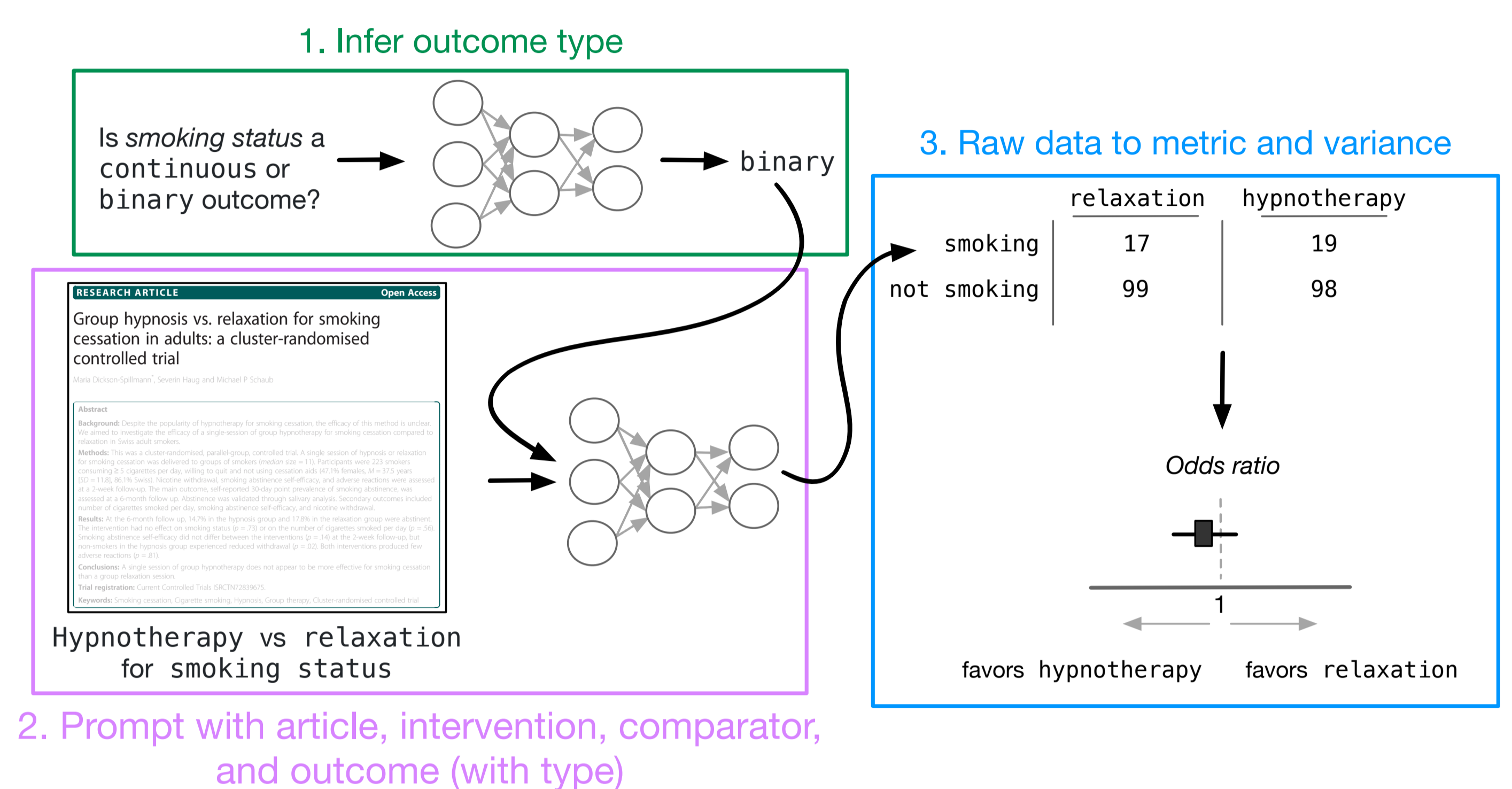
Intervention	Comparator	Outcome	Outcome Type	Intervention Events	Intervention Group Size	Comparator Events	Comparator Group Size
Hypnotherapy	Relaxation	Smoking	Binary	19	116	17	117

Metric	Dev	Test	Total
# PMC Articles	10	110	120
# Prompts (ICOs)	43	656	699
# Binary Outcomes	11	172	183
# Continuous Outcomes	32	484	516
% With Enough Data for Point Estimates	62.79	58.84	59.08
Mean Articles Tokens	3331	3603	3581

## SUMMARY

- **Annotated dataset** for the task of extracting numerical clinical findings for conducting meta-analysis
- **Evaluation** of 8 modern LLMs using the annotated dataset
- **End-to-end case study** of a fully automated meta-analysis
- Binary outcomes extraction: LLMs with large input context windows (e.x. GPT-4) outperform smaller, open-source models
- Continuous outcomes extraction: LLMs perform poorly (below 50% exact match)

## APPROACH



Evaluated 8 LLMs on predicting outcome type and extracting binary and continuous outcomes in YAML format using a **zero-shot approach**. Python's `statsmodels` package was used to derive point estimates and standard errors for meta-analysis.

## RESULTS

### Part 1: Outcome Type

	GPT-4	GPT-3.5	Alpaca	Mistral	Gemma	OLMo	PMC LLaMA	BioMistral
Accuracy	0.713	0.607	<b>0.739</b>	0.201	0.665	0.290	0.732	0.133
F1 - Binary	<b>0.735</b>	0.680	0.000	0.576	0.590	0.424	0.124	0.275
F1 - Continuous	0.836	0.690	<b>0.851</b>	0.183	0.716	0.079	0.848	0.135
# Unknowns	155	152	1	489	0	5	15	409

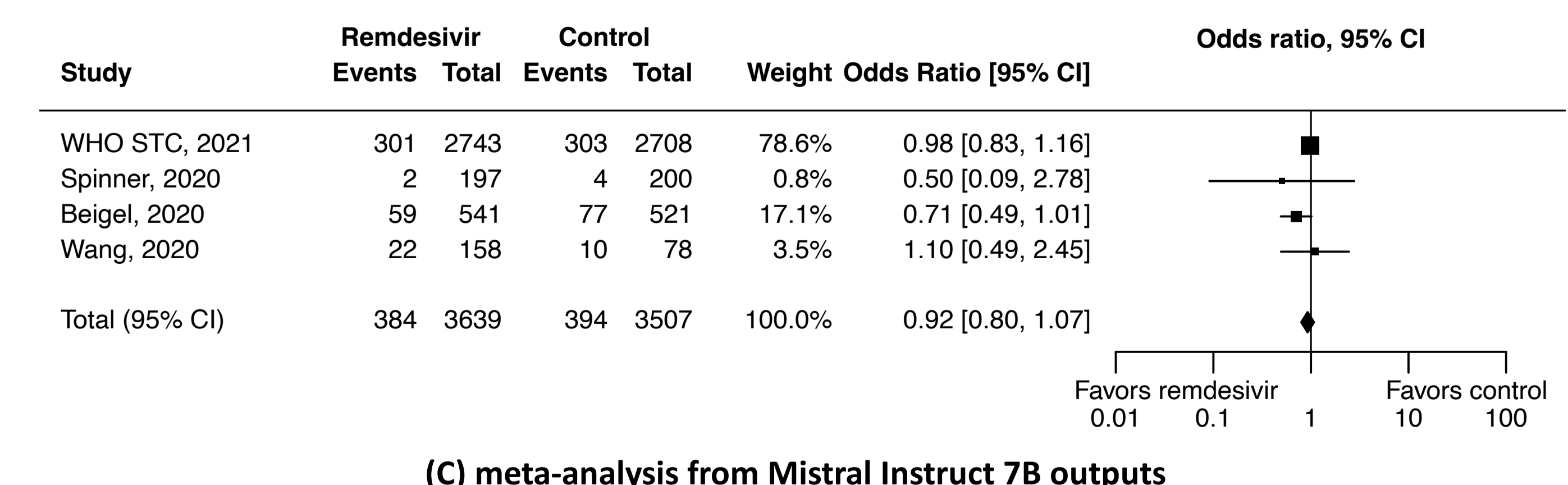
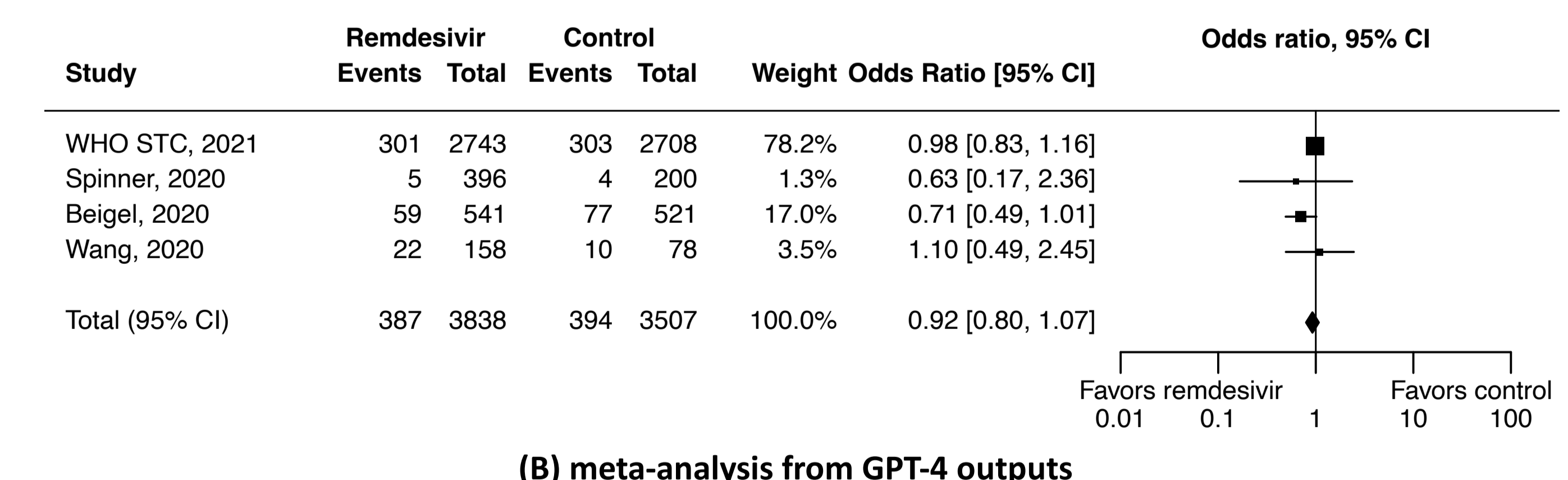
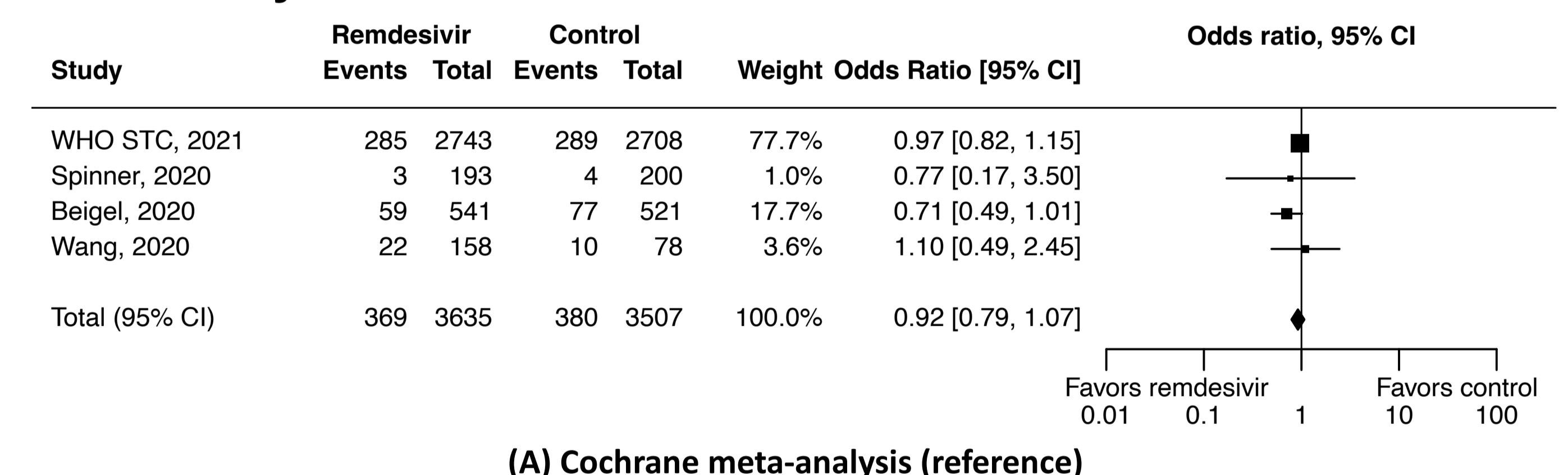
### Part 2a: Binary Outcome Numerical Results Extraction

	GPT-4	GPT-3.5	Alpaca	Mistral	Gemma	OLMo	PMC LLaMA	BioMistral
Exact Match	<b>0.751</b>	0.606	0.334	0.375	0.412	0.311	0.447	0.470
Partial Match	<b>0.912</b>	0.749	0.193	0.708	0.678	0.275	0.216	0.175
MSE	<b>0.101</b>	0.441	0.485	0.657	0.913	1.253	1.523	-
# Unknowns	41	145	490	<b>28</b>	90	319	524	612
% Complete	<b>87.94</b>	61.70	9.22	87.23	58.87	24.11	7.09	0.00

### Part 2b: Continuous Outcome Numerical Results Extraction

	GPT-4	GPT-3.5	Alpaca	Mistral	Gemma	OLMo	PMC LLaMA	BioMistral
Exact Match	<b>0.734</b>	0.641	0.216	0.507	0.534	0.190	0.107	0.087
Partial Match	<b>0.913</b>	0.814	0.470	0.691	0.699	0.408	0.497	0.501
MSE	<b>0.290</b>	0.951	6.257	1.138	3.466	1.738	-	-
# Unknowns	422	437	1169	483	775	1213	1778	1985
% Complete	<b>63.64</b>	62.40	31.82	62.81	40.08	11.98	4.96	0.00

### Case Study: Remdesivir for treatment of COVID-19



Data + Code



This research was partially supported by National Science Foundation (NSF) grants RI-2211954 and IIS-1750978, and by the National Institutes of Health (NIH) under the National Library of Medicine (NLM) grant 2R01LM012086.