

Jiafan Chen, MS¹, Rahul Mudumba, MHS¹, Joshua Morriss, PhD², William V. Padula, PhD^{1,3}

²Ziplitics, Inc., Midlothian, VA USA

³Leonard D. Schaeffer Center for Health Policy and Economics, Los Angeles, CA, USA

We aimed to measure the cost-effectiveness of vaccines for infectious diseases in the U.S. While vaccines are generally cost-effective, their economic value from U.S. healthcare sector and societal perspectives remains under-evaluated¹. Our objective was to use and evaluate the performance of the Literature Review Network (LRN), an explainable AI (XAI), to automate a health outcomes systematic literature review (SLR) and meta-analysis on this topic.

Search Strategy: This SLR was conducted via PubMed using open-access, U.S.-based cost-effectiveness analyses; vaccines were the primary intervention, from a societal/healthcare sector perspective.

Table 1. Search strategy used by LRN. Excluded studies served as a “ground-truth” negative data set.

LRN (v1.5) classified studies using proprietary generative-discriminative models, a word embedding model incorporating the UMLS Metathesaurus, and a meta-heuristic wrapper². Leveraging a large language model (LLM) based on GPT-4-turbo and retrieval-augmented generation (RAG), LRN-included reports were utilized for SLR writing³.

LRN (v1.5) underwent training via reinforcement learning with human feedback (RLHF) over 3 iterations. A final LRN model was deployed to label the entire corpus. During an iteration, the model incorporated feedback from 21 abstracts screened by 2 human researchers, as either INCLUDE or EXCLUDE, and adjustments to the rules. Explainability metrics described concepts LRN used in its decision-making processes. Final included studies were determined by the researchers; a senior health economist broke ties.

Table 2. User-defined parameters across RLHF iterations. “Iteration added” = decisions by researchers during screening.

Post-SLR, LRN (v2.0) was employed to perform a meta-analysis. Using a fixed-effects inverse-variance-weighted model with Fieller's theorem for 95% confidence intervals⁴, the analysis assessed the cost-effectiveness of vaccines through ICERs, with a WTP threshold of \$150,000 (USD/QALY). Data were extracted via LRN (v2.0), with structured tables supported by an LLM based on GPT-4-turbo.

```
graph TD; A[Records identified from:  
Databases (n = 850)  
Registers (n = 0)] --> B[Records screened  
(n = 317)]; A --> C[Records removed before screening:  
Duplicate records (n = 0)  
Records marked as ineligible by LRN too  
(n = 525)  
Records removed for other reasons  
(n = 8)]; B --> D[Reports sought for retrieval  
(n = 300)]; B --> E[Records excluded  
(n = 17)]; D --> F[Reports assessed for eligibility  
(n = 300)]; D --> G[Reports not retrieved  
(n = 0)]; F --> H[New studies included in review  
(n = 154)  
Reports of new included studies  
(n = 154)]; F --> I[Reports excluded:  
LRN Excluded (n = 146)  
SME Excluded with Reason (n = 0)  
Other Reason (n = 0)];
```

Identification of new studies via databases and registers

Records identified from:
Databases (n = 850)
Registers (n = 0)

Records removed before screening:
Duplicate records (n = 0)
Records marked as ineligible by LRN too
(n = 525)
Records removed for other reasons
(n = 8)

Records screened
(n = 317)

Records excluded
(n = 17)

Reports sought for retrieval
(n = 300)

Reports not retrieved
(n = 0)

Reports assessed for eligibility
(n = 300)

Reports excluded:
LRN Excluded (n = 146)
SME Excluded with Reason (n = 0)
Other Reason (n = 0)

New studies included in review
(n = 154)
Reports of new included studies
(n = 154)

Figure 1. Auto-populated PRISMA 2020 flow diagram. After 3 iterations, researchers excluded 17 records and LRN excluded 146 full-text reports. Studies deemed ineligible by LRN (n = 525) were those that met the exclusion criteria. Records removed for other reasons (n = 8) were linguistically incompatible with the NLP processing capabilities of LRN v1.5 (studies written in Chinese or Russian).



Table 3. Productivity metrics for human labor versus computational time via LRN systematic literature review. All time reported in Coordinated Universal Time (UTC), reported in 24-hour time.

Class	Recall	Precision	F-score
	Iteration 1		
INCLUDE	50.00%	56.67%	53.13%
EXCLUDE	41.38%	48.00%	44.44%
	Iteration 2		
INCLUDE	89.19%	63.46%	74.16%
EXCLUDE	17.39%	50.00%	25.81%
	Iteration 3		
INCLUDE	73.53%	64.10%	68.49%
EXCLUDE	46.15%	57.14%	51.06%

Table 4. Performance metrics for XAI classifications. RLHF iteration 3, or the “highest performing iteration” achieved a Cohen’s kappa of 0.2014; highest INCLUDE label performance was in iteration 2.

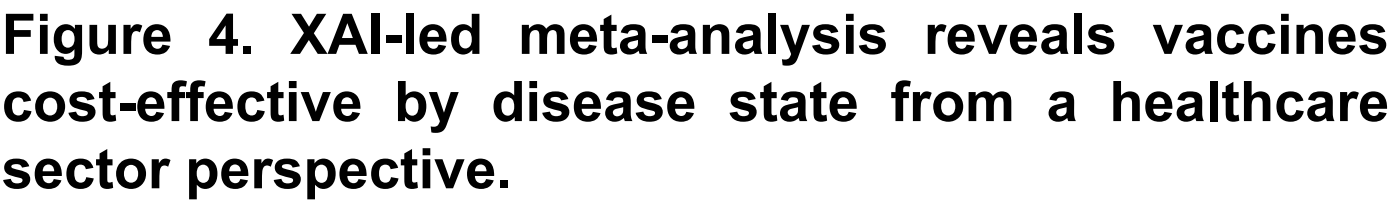


Table 5. Pooled ICERs across CEAs by disease state. HPV = human papillomavirus, HZ = herpes zoster, Hep = hepatitis a/b, MMR = measles, mumps, and rubella, MenB = meningococcal serogroup B, PCV = pneumococcal conjugate vaccine, RZV = recombinant zoster vaccine, Tdap = tetanus, diphtheria, and pertussis.

Our XAI-led SLR and meta-analysis found vaccines to be highly cost-effective from a U.S. healthcare sector perspective, particularly for high-burden diseases such as COVID-19, HIV and HPV at standard WTP thresholds; this supports their preferential status on formularies for health and cost benefits. LRN achieved high performance in 240 minutes. In 300 minutes, the final LRN model screened 850 studies, selecting 154; 777 confirmed by human reviewers, and 45 studies retrieved for final analysis. This study demonstrates that explainable AI, like LRN, is an effective tool for advancing population health.

1. Leidner AJ, Murthy N, Chesson HW, et al. Cost-effectiveness of adult vaccinations: A systematic review. *Vaccine*. 2019;37(2):226-234. doi:10.1016/j.vaccine.2018.11.056
2. Bodenreider, O. (2004). The Unified Medical Language System (UMLS): Integrating biomedical terminology. *Nucleic Acids Research*, 32(Database issue), D267–D270. <https://doi.org/10.1093/nar/gsk061>
3. OpenAI, Achiam, J., Adler, S., ... Zoph, B., ... 2020. [GPT-4 Technical Report \(arXiv:2005.05857\)](https://arxiv.org/abs/2005.05857). arXiv: <https://arxiv.org/abs/2005.05857>
4. Pyle, J., Glick HA, Willis R, et al. Confidence intervals for log odds-ratios as a function of various parameters: a comparison of four methods. *Health Care*. 1997; 8(3):243-247. doi:10.1002/hlth.1098.1050190206.03343:arXiv:2699.0.1.v2