

Sybil-Resistant Service Discovery for Agent Economies

David Shi Kevin Joo
Operator Labs Operator Labs

Oct 31st, 2025

Abstract

x402 enables Hypertext Transfer Protocol (HTTP) services like application programming interfaces (APIs), data feeds, and inference providers to accept cryptocurrency payments for access. As agents increasingly consume these services, discovery becomes critical: which swap interface should an agent trust? Which data provider is the most reliable? We introduce TraceRank, a reputation-weighted ranking algorithm where payment transactions serve as endorsements. TraceRank seeds addresses with precomputed reputation metrics and propagates reputation through payment flows weighted by transaction value and temporal recency. Applied to x402’s payment graph, this surfaces services preferred by high-reputation users rather than those with high transaction volume. Our system combines TraceRank with semantic search to respond to natural language queries with high quality results. We argue that reputation propagation resists Sybil attacks by making spam services with many low-reputation payers rank below legitimate services with few high-reputation payers. Ultimately, we aim to construct a search method for x402 enabled services that avoids infrastructure bias and has better performance than purely volume based or semantic methods.

1 Introduction

Autonomous AI systems need machine-native payments to operate without humans in the loop. Legacy rails enforce subscriptions, delayed settlement, chargebacks, and manual invoicing, which block the real-time procurement of context, workflows, and compute. Browser usage reduce some friction but still assume human UX and credit cards. The x402 protocol [1] closes the payments gap with onchain HTTP, so agents can pay per request and receive responses atomically (e.g., in USDC [2]). Our complementary goal is discovery: once agents can pay for any service, how do they decide *which* to call?

Our central insight is to treat each payment as an endorsement whose strength depends on payer reputation, turning discovery into reputation-weighted ranking over the payment

graph. We contribute: (i) a concise failure analysis of count/volume baselines and unseeded PageRank [3]; (ii) TraceRank, which seeds reputation and propagates it along value- and time-weighted flows (resolving the limits of manual TrustRank style ranking [4]); and (iii) a minimal system that combines TraceRank with vector retrieval to return reliable services for agent queries.

2 Background: x402 Payment-Gated Services

x402 standardizes payment-gated HTTP endpoints: a client pays a service address, the service responds, and the payment is recorded onchain. Each payment forms a directed edge from payer to service with value and timestamp. This transparency enables discovery via revealed preference.

Concretely, an initial request to a paid endpoint receives an HTTP 402 (Payment Required) describing terms; the client replays the request with a signed payment payload (e.g., an ‘X-PAYMENT’ header). The server or an optional facilitator verifies the payment, triggers on-chain settlement, and returns a 200 OK with the resource and an ‘X-PAYMENT-RESPONSE’ containing settlement details. The protocol is stateless, HTTP-native, and chain-agnostic via facilitators, making per-request settlement practical for autonomous agents [1].

Perhaps the most compelling characteristic of the x402 protocol is the simplicity of integrating the client and server packages and its subsequent permissionless nature. Anyone can build a server, a facilitator, or a client, and with the increasing capability of large language models (LLMs), any agent can bootstrap the aforementioned with ease.

3 Why Existing Approaches Fail

Counting payments invites Sybil spam and elevates infrastructure. Ranking by total volume rewards whales and enables wash trading. Both approaches over-index on quantity while ignoring who is paying and whether usage legitimate.

Unseeded PageRank [3] on payment graphs promotes compulsory contracts (stablecoin issuers, automated market makers (AMMs), routers) rather than services, because transactions are protocol-mandated rather than editorial judgments. Manual seed lists (e.g., TrustRank [4]) do not generalize in pseudonymous settings. None of these methods tie endorsements to who paid, how much, and when.

3.1 The Spam Service Problem

Consider two x402 services competing for discovery. *Service A* (spam) advertises “Send \$1, receive 1M airdrop” and attracts 10,000 fresh wallets (airdrop farmers, bots), totaling 10,000

payments and \$10,000 volume. *Service B* (legitimate) is a high-quality, specific background check service used by 50 sophisticated traders and protocols, totaling 50 payments and \$5,000 volume.

Under naive ranking, *Service A* wins on transaction count ($10,000 \gg 50$), on volume ($10K \gg 5K$), and on unseeded PageRank (more inbound edges). Popularity among low-reputation users outweighs endorsement by high-reputation users, inverting the discovery objective: agents searching for reliability see spam.

TraceRank weights each payment by payer reputation. If *Service A*'s payers have near-zero seeds (fresh wallets), their collective endorsement contributes negligible reputation. If *Service B*'s payers have high seeds (proven traders), each payment carries substantial weight. The outcome reverses: *Service B* ranks higher despite lower raw counts and volume.

4 TraceRank: Payments as Endorsements

TraceRank operationalizes the payments-as-endorsements insight while resisting infrastructure bias and Sybil attacks. We precompute per-address reputation scores from external signals, then propagate those scores through economically and temporally weighted payment flows. Rather than counting transactions uniformly, each payment is weighted by the payer's seed reputation, value, and recency.

Let \mathcal{V} denote addresses (users and services) and \mathcal{E} denote directed payment edges. For each address $i \in \mathcal{V}$, assign a seed score s_i from external data; addresses without seed data receive $s_i = 0$. Over a chosen observation window, define aggregated flow

$$F_{j \rightarrow i} = \sum_{e \in \mathcal{E}_{j \rightarrow i}} \left[\log(1 + \text{value}_{\text{USD}}(e)) \cdot e^{-\lambda \text{age}_{\text{days}}(e)} \right], \quad (1)$$

where \log is natural and λ has units of day^{-1} (ages measured in days). Let $S_i = \sum_k F_{k \rightarrow i}$ and define the normalized incoming-flow matrix W by

$$w_{ji} = \begin{cases} \frac{F_{j \rightarrow i}}{S_i}, & S_i > 0, \\ 0, & S_i = 0 \text{ (sinks)}, \end{cases} \quad (2)$$

so that columns with positive inbound flow are exactly column-stochastic. TraceRank iterates

$$\mathbf{r}^{(t+1)} = \mathbf{s} + \alpha W^\top \mathbf{r}^{(t)}, \quad \alpha \in (0, 1), \quad (3)$$

which converges to the unique fixed point

$$\mathbf{r} = (I - \alpha W^\top)^{-1} \mathbf{s}. \quad (4)$$

An immediate consequence is Sybil resistance: if a service receives payments from N addresses with zero seed scores, it accumulates zero propagated reputation regardless of N . Bot wallets contribute no signal. Conversely, a single payment from a high-seed payer propagates meaningful reputation. This asymmetry makes fake popularity economically unappealing. The score r_i reflects both direct seed reputation and propagated endorsements along recent, value-weighted flows.

TraceRank is agnostic to seed provenance. Seeds may combine trading performance, decentralized social signals (e.g., Farcaster [6]), protocol contributions, labeled entities (decentralized autonomous organizations (DAOs), verified protocols, funds), and agent attestations from ERC-8004 registries (identity, reputation, validation) [5]. Log scaling, temporal decay, and normalization curb whales and wash trading, surface recency, and prevent infrastructure accumulation.

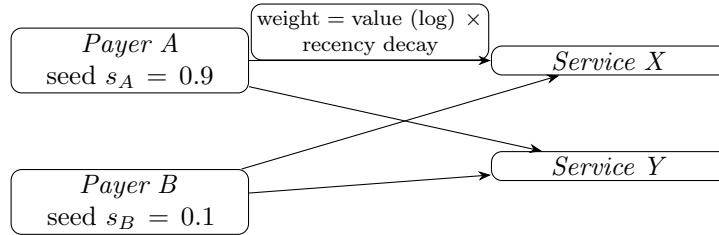


Figure 1: Reputation propagation: high-seed payers diffuse reputation along value- and time-weighted payment flows to services.

5 System Architecture for x402 Service Discovery

We precompute a single TraceRank score per service from seeded accounts with value and time-weighted payment flows. Profiles are natural language descriptions of what the service does, and the vector embeddings of these profiles are precomputed. We use dense vector retrieval (optionally hybrid sparse + dense) to retrieve top- K by cosine similarity and fuse with TraceRank via a simple multiplicative score:

$$\text{score}(A, q) = \cos(\mathbf{q}, \mathbf{p}_A) \times \text{TraceRank}(A). \quad (5)$$

This ranking returns the most semantically relevant services among those preferred by reputable payers. Agents can refine results by rephrasing or expanding the query and rerunning retrieval.

In an example PostgreSQL implementation with pgvector, store a per-service ‘tracerank’ column and embeddings in one table, add vector and btree indexes, and compute the final score directly in SQL (e.g., ‘ORDER BY cosine(embedding, :q) * tracerank DESC’) with optional ‘WHERE’ filters for chain, time window, or tags.

Counterfactual techniques include: (i) semantic-only, (ii) TraceRank-only, and (iii) volume

and count oriented query techniques. Future work will demonstrate how the combined TraceRank and vector similarity technique excels at retrieving the highest quality services.

6 Conclusion

Quality emerges from who pays, not just how much. TraceRank propagates precomputed reputation through value and time-weighted payment flows, producing a single score per service that resists Sybil spam, whale dominance, and infrastructure bias.

Combined with semantic retrieval of service profiles, a simple multiplicative fusion yields a fast, agent-ready ranking. As agent economies scale, fast and high-quality service discovery becomes critical infrastructure. TraceRank shows that reputation can emerge from payment patterns and precomputed address related social reputation, enabling agents to bootstrap trust in decentralized marketplaces without privileged curators or identity systems.

Acknowledgements

We thank SM Mesbahul Islam for his contributions to the reference implementation and valuable engineering support.

References

- [1] E. Reppel, R. Caspers, K. Leffew, D. Organ, D. Kim, and N. Dalal, “x402: An open standard for internet-native payments,” Coinbase Developer Platform whitepaper, May 6, 2025.
- [2] Centre Consortium, “USD Coin (USDC),” documentation, accessed 2025.
- [3] S. Brin and L. Page, “The anatomy of a large-scale hypertextual web search engine,” *Computer Networks and ISDN Systems*, vol. 30, no. 1–7, pp. 107–117, 1998.
- [4] Z. Gyöngyi, H. Garcia-Molina, and J. Pedersen, “Combating web spam with TrustRank,” in *VLDB*, 2004.
- [5] M. De Rossi, D. Crapis, J. Ellis, and E. Reppel, “ERC-8004: Trustless Agents [DRAFT],” Ethereum Improvement Proposals, no. 8004, Aug. 2025. Available: <https://eips.ethereum.org/EIPS/eip-8004>.
- [6] V. Srinivasan, D. Romero, et al., “Farcaster: A Decentralized Social Network,” documentation, accessed 2025.