

# Deep Reinforcement Learning-based Rebalancing Policies for Profit Maximization of Relay Nodes in Payment Channel Networks

**Nikolaos Papadis and Leandros Tassiulas**

Department of Electrical Engineering & Institute for Network Science  
Yale University

{nikolaos.papadis, leandros.tassiulas}@yale.edu

October 13, 2022

## **Abstract**

Payment channel networks (PCNs) are a layer-2 blockchain scalability solution, with its main entity, the payment channel, enabling transactions between pairs of nodes “off-chain,” thus reducing the burden on the layer-1 network. Nodes with multiple channels can serve as relays for multihop payments over a path of channels: they relay payments of others by providing the liquidity of their channels, in exchange for part of the amount withheld as a fee. Relay nodes might after a while end up with one or more unbalanced channels, and thus need to trigger a rebalancing operation. In this paper, we study how a relay node can maximize its profits from fees by using the rebalancing method of submarine swaps. We introduce a stochastic model to capture the dynamics of a relay node observing random transaction arrivals and performing occasional rebalancing operations, and express the system evolution as a Markov Decision Process. We formulate the problem of the maximization of the node’s fortune over time over all rebalancing policies, and approximate the optimal solution by designing a Deep Reinforcement Learning (DRL)-based rebalancing policy. We build a discrete event simulator of the system and use it to demonstrate the DRL policy’s superior performance under most conditions by conducting a comparative study of different policies and parameterizations. In all, our approach aims to be the first to introduce DRL for network optimization in the complex world of PCNs.

# 1 Introduction

Blockchain technology enables trusted interactions between untrusted parties, with financial applications like Bitcoin, and beyond, but with also known scalability issues [1]. Payment channels are a layer-2 development towards avoiding the long confirmation times and high costs of the layer-1 network: they enable nodes that want to transact quickly, cheaply and privately to do so by depositing some balances to open a payment channel between themselves, and then trustlessly shift the same total balance between the two sides without broadcasting their transactions and burdening the network. Connected channels create a Payment Channel Network (PCN), via which two nodes not sharing a channel can still pay one another via a sequence of existing channels. Intermediate nodes in the PCN function as relays: they forward the payment along its path and collect relay fees in return. As transactions flow through the PCN, some channels get depleted, causing incoming transactions to fail because of insufficient liquidity on their path. Thus, the need for channel rebalancing arises.

In this paper, we study the rebalancing mechanism of submarine swaps, which allows a blockchain node to exchange some funds from on- to off-chain and vice versa. Since a swap involves an on-chain transaction, it takes some time to complete. Taking this into account, we formulate the following optimal rebalancing problem as a Markov Decision Process (MDP): For a node relaying traffic across multiple channels, determine an optimal rebalancing strategy over time (i.e. when and how much to rebalance) as a function of the transaction arrival rates observed from an unknown distribution and the confirmation time of an on-chain transaction, so that the node can keep its channels liquid and its profit from relay fees can be maximized.

More specifically, our contributions are the following:

- We develop a stochastic model that captures the dynamics of a relay node in an environment of two timescales: a continuous one for random discrete transaction arrivals in both directions from distributions unknown to the node, and a discrete one for dispatching rebalancing operations.
- We express the system evolution in our model as an MDP with continuous state and action spaces and time-dependent constraints on the actions, and formulate the problem of relay node profit maximization.
- We approximate the optimal policy of the MDP using Deep Reinforcement Learning (DRL) by appropriately engineering the states, actions and rewards and tuning a version of the Soft Actor-Critic algorithm.
- We develop a realistic discrete event simulator of the system, and use it to evaluate the performance of the learning-based as well as other heuristic rebalancing policies under various transaction arrival conditions and demonstrate the superiority of our policy in a range of regimes.

In summary, our paper is the first to formally study the submarine swap rebalancing mechanism and to introduce a DRL-based method for channel rebalancing in particular, and for PCN liquidity management in general.

The remainder of the paper is organized as follows. In Sec. 2 we introduce the operation of payment channels and relay nodes, explain the need for rebalancing, and introduce the submarine swap rebalancing mechanism. In Sec. 3 we describe our stochastic model of a relay node with two payment channels and write the profit maximization problem using an MDP. In Sec. 4 we present heuristic policies as well as design a DRL-based algorithm for an approximately optimal solution to the problem, and in Sec. 5 we describe the experimental setup and results. In Sec. 6 and 7 we discuss future directions and some related work. Finally, Sec. 8 concludes the paper.

## 2 Background

### 2.1 Payment channel operation

A payment channel (Fig. 1) is created between two nodes  $N_1$  and  $N_2$  after they deposit some capital to a channel-opening on-chain transaction. After this transaction is confirmed, the nodes can transact completely *off-chain* (i.e. in the channel) without broadcasting their interactions to the layer-1 network, and without the risk of losing funds thanks to a cryptographic safety mechanism. The sum of their two balances in the channel remains constant and is called the channel capacity. A transaction of amount  $\alpha$  from  $N_1$  to  $N_2$  will succeed if the balance of  $N_1$  at that moment suffices to cover it. In this case, the balance of  $N_1$  is reduced by  $\alpha$  and the balance of  $N_2$  is increased by  $\alpha$ .

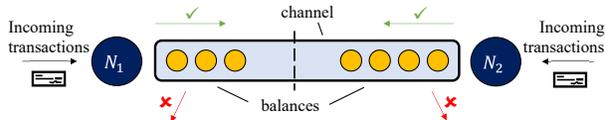


Figure 1: A payment channel between nodes  $N_1$  and  $N_2$  and current balances of 3 and 4

### 2.2 The role of relay nodes

As pairs of nodes create channels, a payment channel network (Fig. 2) is formed, over which multihop payments are possible: Consider a transaction of amount 5 from  $N_1$  to  $N_3$  via  $N_2$ . Note that the amount 5 includes the fees it will have to pay on its way, e.g. 1% at each intermediate node. In the  $N_1N_2$  channel,  $N_1$ 's local balance is reduced by 5 and  $N_2$ 's local balance is increased by 5. In the  $N_2N_3$  channel,  $N_2$ 's local balance is reduced by  $5 - \text{fees} = 4.95$  and  $N_3$ 's local balance is increased by 4.95.  $N_2$ 's total capital in all its channels before the transaction was  $2 + 1 + 7 = 10$ , while after it is  $7 + 1 + 2.05 = 10.05$ , so  $N_2$  made a profit of 0.05 by acting as a relay. If one of the outgoing balances did not suffice, then the transaction would fail end-to-end, thanks to a cryptographic smart contract mechanism called Hashed Time-Lock Contract (HTLC). The role of relay nodes is funda-

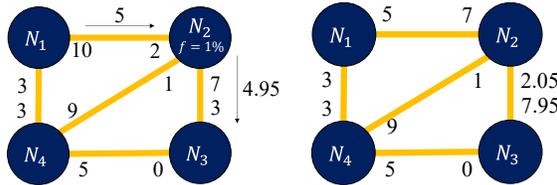


Figure 2: Processing of a transaction in a payment channel network: before (left) and after (right)

mental for the continuous operation of a PCN. The most prominent PCN currently is the Lightning Network [2] built on top of Bitcoin. More details on PCN operation can be found in [3, 4].

## 2.3 The need for rebalancing

Depending on the demand a payment channel is facing in its two directions, funds might accumulate on one side and deplete on the other. This might happen due to asymmetric demand inside single channels, the random nature of arrivals causing temporary depletions at specific times (e.g. when a large transaction arrives), or even symmetric demand between two endpoints of a multihop path which can cause imbalance due to fees withheld by intermediate nodes (an example of this latter more subtle case is given in Appendix A). The resulting imbalance is undesirable, as it leads to transaction failures and loss of profit from relay fees. In fact, an entire PCN will stop being operational in finite time without external intervention, creating the need for rebalancing mechanisms.

## 2.4 The submarine swap rebalancing mechanism

In this work, we study submarine swaps, introduced in [5] and used commercially by Boltz<sup>1</sup> and Loop<sup>2</sup>. At a high level, a *submarine swap* works as follows (Fig. 3): Node  $N_1$  owns some funds in its channel with node  $N_2$ , and some funds on-chain. At time  $t_0$ , the channel  $N_1N_2$  is almost depleted on  $N_1$ 's side (balance = 5).  $N_1$  can start a *swap-in* by paying an amount (50) to a Liquidity Service Provider (LSP) - a wealthy node with access to both layers - via an *on-chain* transaction, and the LSP will give this amount back (reduced by a 10% swap fee, so 45) to  $N_1$  *off-chain* via a path that goes through  $N_2$ . The final amount that is added at  $N_1$  (and subtracted at  $N_2$ ) is  $45 - \epsilon$  due to the relay fees spent on its way from the LSP. Thus, at time  $t_1$  the channel will be almost perfectly balanced. The reverse procedure is also possible (a *reverse submarine swap* or *swap-out*) in order for a node to offload funds from its channel, by paying the server off-chain and receiving funds on-chain.<sup>3</sup>

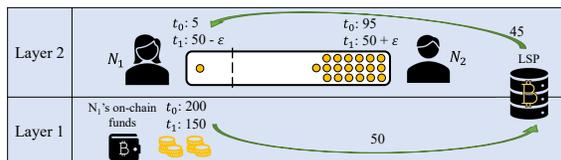


Figure 3: A submarine swap (swap-in)

A sketch of the technical protocol followed during a successful swap-in (a swap-out is similar) is shown in Fig. 4 and subsequently modeled in Sec. 3.1.3. First, a node-client initiates the swap by generating a hash preimage, creating an invoice of the desired swap amount  $r$  tied to this hash and with a certain expiration time  $T_{exp}$ , and sends it to an LSP that is willing to make the exchange.

The LSP then quotes what it wants to be paid on-chain in exchange for paying the client's invoice off-chain, say  $\alpha + F_{swap}(\alpha)$ , where  $F_{swap}(\alpha)$  is the LSP's swap service fee. If the client accepts the exchange rate, it creates a conditional on-chain payment of amount  $\alpha + F_{swap}(\alpha)$

<sup>1</sup><https://boltz.exchange>

<sup>2</sup><https://lightning.engineering/loop>

<sup>3</sup>In both cases, the layer-2 channel balances are altered with the help of an on-chain (layer-1, one layer below the channel) transaction, hence the characterization "submarine."

to the LSP based on an HTLC with the same preimage as before and broadcasts the payment to the blockchain network. The payment can only be redeemed if the LSP knows the preimage, and the client will only reveal the preimage once it has received the LSP’s funds on-chain. Thus, the LSP pays the off-chain invoice. This forces the client to reveal the preimage, and now the LSP can redeem the on-chain funds and the swap is complete. The entire process happens trustlessly thanks to the HTLC mechanics. More technical details can be found in [6, 7].

There is an important tradeoff the node has to make, which is to strike a balance in terms of how often and how much it should rebalance: it can choose to not rebalance a lot to avoid paying swap fees, but then it forfeits profits from relay fees of transactions dropped due to imbalance. On the other hand, it can choose to rebalance a lot so as not to drop any transaction, but then entails high rebalancing fee costs. This observation motivates us to study the problem of demand-aware and timely dispatching of swaps of the right amount by a node aiming to maximize its total fortune, which is presented in the next section.

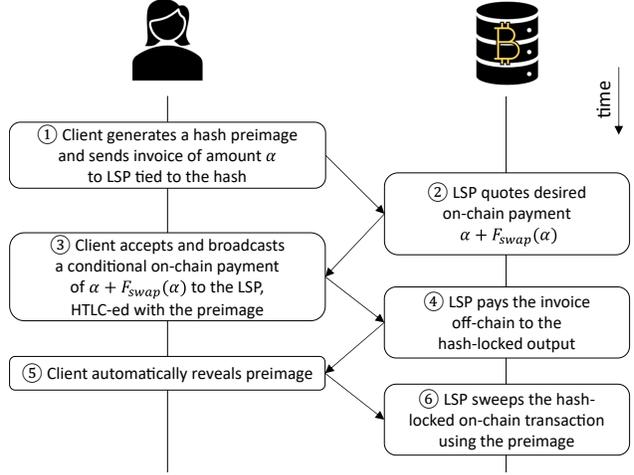


Figure 4: A swap-in step-by-step.

### 3 Problem formulation

#### 3.1 System evolution

In this section, we introduce a stochastic model of a PCN relay node  $N$  that has two channels, one with node  $L$  and one with node  $R$ , and wishes to maximize its profits from relaying payments from  $L$  to  $R$  and vice versa (Fig. 5). Define  $b_{LN}(\tau)$ ,  $b_{NL}(\tau)$ ,  $b_{NR}(\tau)$ ,  $b_{RN}(\tau)$  to be the balances of the channels and  $B_N(\tau)$  to be the on-chain amount of  $N$  at time  $\tau$ . Let  $C_n$  be the total capacity of the channel  $Nn$ ,  $n \in \mathcal{N} \triangleq \{L, R\}$ . Events happen at two timescales: a continuous one for arriving transactions, and a discrete one for times when the node is allowed to rebalance.

##### 3.1.1 The transaction timescale

Transactions arrive as a marked point process and are characterized by their direction ( $L$ -to- $R$  or  $R$ -to- $L$ ), time of arrival and amount<sup>4</sup>. At each moment in continuous time (denoted by  $\tau$ ), (at most) one transaction arrives in the system. All transactions are admitted but some fail due to insufficient balances.

<sup>4</sup>We consider node  $N$  to not be the source or destination of any transactions itself, but rather only to act as a relay.

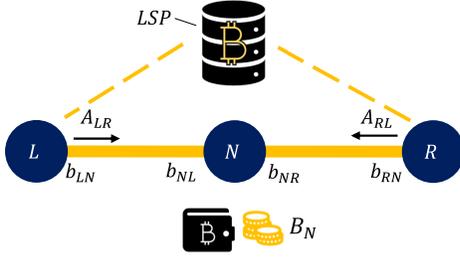


Figure 5: System model

and coming from node  $R$  in the  $R$ -to- $L$  direction at time  $\tau$  respectively, each drawn from a distribution that is fixed but unknown to node  $N$ . An arriving transaction of amount  $A_{LR}(\tau) = \alpha$  is feasible if and only if there is enough balance in the  $L$ -to- $R$  direction in both channels, i.e.  $b_{LN}(\tau) \geq \alpha$  and  $b_{NR} \geq \alpha - f(\alpha)$ , and similarly for the  $R$ -to- $L$  direction. The successfully admitted and processed amounts by node  $N$  at time  $\tau$  are<sup>5</sup>:

$$S_{LR}(\tau) = \begin{cases} A_{LR}(\tau) & , \text{ if } A_{LR}(\tau) \leq b_{LN}(\tau) \text{ and } A_{LR}(\tau) - f(A_{LR}(\tau)) \leq b_{NR}(\tau) \\ 0 & , \text{ otherwise} \end{cases} \quad (1)$$

$$S_{RL}(\tau) = \begin{cases} A_{RL}(\tau) & , \text{ if } A_{RL}(\tau) \leq b_{RN}(\tau) \text{ and } A_{RL}(\tau) - f(A_{RL}(\tau)) \leq b_{NL}(\tau) \\ 0 & , \text{ otherwise} \end{cases} \quad (2)$$

Then the profit of node  $N$  at time  $\tau$  is  $f(S_{LR}(\tau)) + f(S_{RL}(\tau))$ , and the lost fees (from transaction that potentially failed to process) are  $f(A_{LR}(\tau) - S_{LR}(\tau)) + f(A_{RL}(\tau) - S_{RL}(\tau))$ . The balance processes at time  $\tau$  evolve as follows:

$$b_{LN}(\tau) \rightarrow b_{LN}(\tau) + (S_{RL}(\tau) - f(S_{RL}(\tau))) - S_{LR}(\tau) \quad (3)$$

$$b_{NL}(\tau) \rightarrow b_{NL}(\tau) + S_{LR}(\tau) - (S_{RL}(\tau) - f(S_{RL}(\tau))) \quad (4)$$

$$b_{NR}(\tau) \rightarrow b_{NR}(\tau) + S_{RL}(\tau) - (S_{LR}(\tau) - f(S_{LR}(\tau))) \quad (5)$$

$$b_{RN}(\tau) \rightarrow b_{RN}(\tau) + (S_{LR}(\tau) - f(S_{LR}(\tau))) - S_{RL}(\tau) \quad (6)$$

The on-chain amount  $B_N(\tau)$  is not affected by the processing of off-chain transactions. The system model is shown in Fig. 5.

### 3.1.2 The rebalancing decision (control) timescale

The evolution of the system can be controlled by node  $N$  using “submarine swap” rebalancing operations. Rebalancing may start at times  $t_i = i \cdot T_{check}$ ,  $i = 0, 1, \dots$ , and takes a (fixed)

<sup>5</sup>Since in the sequel we focus on the discrete and sparse time scale of the periodic times at which the node rebalances, we make the fair assumption (as e.g. in [9]) that off-chain transactions are processed instantaneously across their entire path and do not fail in their subsequent steps after they cross the two channels (if a transaction were to fail outside the two channels, it can be viewed as of zero value by the system).

time  $T_{conf}$  to complete (on average 10 minutes for Bitcoin)<sup>6</sup>. We consider the case where  $T_{check} \geq T_{conf}$  (to avoid having concurrent rebalancing operations in the same channel that could be combined into one).

The system state is defined only for the discrete rebalancing decision timescale as the collection of the off- and on-chain balances:

$$S(t_i) = (b_{LN}(t_i), b_{NL}(t_i), b_{NR}(t_i), b_{RN}(t_i), B_N(t_i)) \quad (7)$$

At each time  $t_i$ , node  $N$  can decide to commence a swap-in or a swap-out in each channel. Call the respective amounts  $r_L^{in}(t_i), r_L^{out}(t_i), r_R^{in}(t_i), r_R^{out}(t_i)$ . At any time  $t_i$ , in a given channel, either a swap-in or a swap-out or nothing will be commenced, but not both a swap-in and a swap-out<sup>7</sup>.

Let  $F_{swap}^{in}(\alpha)$  and  $F_{swap}^{out}(\alpha)$  be the swap fees that the LSP charges for an amount  $\alpha$  for a swap-in and a swap-out respectively, where  $F_{swap}^{in}(\cdot)$  and  $F_{swap}^{out}(\cdot)$  are any functions with  $F_{swap}(0) = 0$ . For ease of exposition, we let all types of fees the node will have to pay (relay fees for the off-chain part, on-chain miner fees, server fees) be both part of the above swap fees, and be the same for swap-in and swap-out when a net amount  $r_{net}$  is transferred from on- to off-chain or vice versa:  $F_{swap}^{in}(r_{net}) = F_{swap}^{out}(r_{net}) = F_{swap}(r_{net}) \triangleq r_{net}F + M$ , where the proportional part  $F$  includes the server fee and off-chain relay fees, and  $M$  includes the miner fee and potential base fees.

Note that the semantics of the swap amounts  $r$  are such that they represent the amount that will move *in the channel* (and not necessarily the net change in the node's fortune). As a result of this convention and based on the swap operation as described in section 3.1.3, the amount  $r^{in}$  of a swap-in coincides with the net amount  $r_{net}^{in}$  by which the node's fortune decreases (as  $r^{in}$  does not include the swap fee), while the amount  $r^{out}$  of a swap-out includes the swap fee and the net amount by which the node's fortune decreases is  $r_{net}^{out} = \phi^{-1}(r^{out})$ , where  $\phi(r_{net}) \triangleq r_{net} + F_{swap}(r_{net})$ , and  $\phi^{-1}$  the generalized inverse function of  $\phi(\cdot)$  (it always exists:  $\phi^{-1}(y) = \min\{x \in \mathbb{N} : \phi(x) = y\}$ ). For our  $F_{swap}(\cdot)$  it is  $\phi(r_{net}) = r_{net}(1 + F) + M$  for  $r_{net} > 0$ ,  $\phi(0) = 0$ , so  $\phi^{-1}(y) = (y - M)/(1 + F)$  for  $y > 0$  and  $\phi^{-1}(0) = 0$ .

### 3.1.3 A submarine swap step-by-step

We now describe how a rebalancing operation on the  $Nn$  channel is affecting the system state. First, a swap-in of amount  $r_n^{in}$  initiated by node  $N$  to refill  $N$ 's local balance in the  $Nn$  channel:

- At time  $t_i$ , node  $N$  locks the net rebalancing amount plus fees and subtracts it from its on-chain funds:  $B_N \rightarrow B_N - (r_n^{in} + F_{swap}^{in}(r_n^{in}))$
- At time  $t_i + T_{conf}$ , the on-chain transaction is confirmed, so the LSP sends a payment of  $r_n^{in}$  to node  $N$  off-chain<sup>8</sup>. The rebalancing payment reaches node  $n$ :

<sup>6</sup>In practice, completion happens when the miners solve the random puzzle and produce the Proof-of-Work for the next block that includes the rebalancing transaction. The time for this to happen fluctuates, though only slightly, so we use a fixed value for the sake of tractability.

<sup>7</sup>The other two nodes ( $L$  and  $R$ ) are considered passive, i.e. they perform no swap operations themselves.

<sup>8</sup>The LSP is a well-connected node owning large amounts of liquidity, so we reasonably assume that it can always find a route from itself to  $N$ , possibly via splitting the amount across multiple paths.

If  $b_{nN} \leq r_n^{in}$  (i.e.  $n$  does not have enough balance to forward it), then the off-chain payment fails. The on-chain funds are unlocked<sup>9</sup> and refunded back to the on-chain amount:  $B_N \rightarrow B_N + (r_n^{in} + F_{swap}^{in}(r_n^{in}))$

Otherwise (if the transaction is feasible),  $n$  forwards the payment to  $N$ :  $b_{nN} \rightarrow b_{nN} - r_n^{in}$  and  $b_{Nn} = b_{Nn} + r_n^{in}$

A swap-out of amount  $r_n^{out}$ , initiated by node  $N$  to offload some of its local balance to the chain, works as follows:

- At time  $t_i$ , node  $N$  locks the net rebalancing amount plus fees and sends it to the LSP via the off-chain network:  $b_{Nn} \rightarrow b_{Nn} - r_n^{out}$ . Note that  $r_n^{out}$  includes the fees.
- At time  $t_i + T_{conf}$ , the on-chain transaction is confirmed, so node  $N$  receives the funds on-chain:  $B_N \rightarrow B_N + \phi^{-1}(r_n^{out})$ , and the funds are also unlocked in the channel and pushed towards the remote balance:  $b_{nN} \rightarrow b_{nN} + r_n^{out}$

### 3.1.4 Rebalancing constraints

Based on the steps just described, the rebalancing operations will succeed if and only if their amounts satisfy the following constraints:

Rebalancing amounts must be non-negative:

$$r_n^{in}(t_i), r_n^{out}(t_i) \geq 0 \text{ for all } i \in \mathbb{N}, n \in \mathcal{N} \quad (8)$$

A swap-in and a swap-out cannot be requested in the same channel at the same time<sup>10</sup>:

$$r_n^{in}(t_i) \cdot r_n^{out}(t_i) = 0 \text{ for all } i \in \mathbb{N}, n \in \mathcal{N} \quad (9)$$

The swap-out amounts (which already include the swap fees) must be greater than the fees themselves:

$$r_n^{out}(t_i) - F_{swap}^{out}(r_n^{out}(t_i)) \geq 0 \text{ for all } i \in \mathbb{N}, n \in \mathcal{N} \quad (10)$$

The respective channel balances must suffice to cover the swap-out amounts (which already include the swap fees):

$$r_n^{out}(t_i) \leq b_{Nn}(t_i) \text{ for all } i \in \mathbb{N}, n \in \mathcal{N} \quad (11)$$

The on-chain balance must suffice to cover the total swap-in amount plus fees:

$$\sum_{n \in \mathcal{N}} (r_n^{in}(t_i) + F_{swap}^{in}(r_n^{in}(t_i))) \leq B_N(t_i) \text{ for all } i \in \mathbb{N} \quad (12)$$

---

<sup>9</sup>In practice, the on-chain funds are unlocked after a time  $T_{exp}$  to prevent malicious clients from requesting many swaps from an LSP and then defaulting. However, since we are concerned with online and cooperative clients with on-chain amounts usually quite larger than the amounts in their channels (and thus than their swaps), and also there is currently a community effort to reduce or even eliminate  $T_{exp}$ , we ignore it.

<sup>10</sup>This fact allows us to express the decision per channel as a single variable taking both positive and negative values, instead of two non-negative variables. We do so in Section 4.2, but we retain two action variables per channel in the present section for the sake of clarity.

### 3.1.5 State evolution equations

Now we are able to write the complete state evolution equations. The amounts added to each balance due to successful transactions during the interval  $(t_i, t_{i+1})$  are

$$d_{NL}^{(t_i, t_{i+1})} \triangleq \int_{\tau \in (t_i, t_{i+1})} \left( S_{LR}(\tau) - (S_{RL}(\tau) - f(S_{RL}(\tau))) \right) d\tau \quad (13)$$

$$d_{NR}^{(t_i, t_{i+1})} \triangleq \int_{\tau \in (t_i, t_{i+1})} \left( S_{RL}(\tau) - (S_{LR}(\tau) - f(S_{LR}(\tau))) \right) d\tau \quad (14)$$

and  $d_{nN}^{(t_i, t_{i+1})} \triangleq -d_{Nn}^{(t_i, t_{i+1})}$ . Then for actions taken subject to the constraints (8)–(12), the state evolves as follows:

$$b_{nN}(t_{i+1}) = b_{nN}(t_i) + d_{nN}^{(t_i, t_{i+1})} - (r_n^{in}(t_i) - z_n(t_i)) + r_n^{out}(t_i) \quad (15)$$

$$b_{Nn}(t_{i+1}) = b_{Nn}(t_i) + d_{Nn}^{(t_i, t_{i+1})} + (r_n^{in}(t_i) - z_n(t_i)) - r_n^{out}(t_i) \quad (16)$$

$$B_N(t_{i+1}) = B_N(t_i) - \sum_{n \in \mathcal{N}} \left( r_n^{in}(t_i) - F_{swap}(r_n^{in}(t_i)) \right) + \sum_{n \in \mathcal{N}} \phi^{-1}(r_n^{out}(t_i)) + \sum_{n \in \mathcal{N}} w_n(t_i) \quad (17)$$

where  $z_n(t_i)$  and  $w_n(t_i)$  are the refunds of the swap-in amount off- and on-chain respectively in case a swap-in operation fails:

$$z_n(t_i) = r_n^{in}(t_i) \mathbb{1}\{b_{nN}(t_i) + d_{nN}^{(t_i, t_i + T_{conf})} < r_n^{in}(t_i)\} \quad (18)$$

$$w_n(t_i) = z_n(t_i) + F_{swap}^{in}(z_n(t_i)) \quad (19)$$

## 3.2 Writing the problem as a Markov Decision Process

The objective function the node wishes to maximize in the real world is its *total fortune both in the channels and on-chain*. The fortune increase due to the action (the 4-tuple)  $r(t_i)$  taken at step  $t_i$  is:

$$D(t_i, r(t_i)) \triangleq \left( \sum_{n \in \mathcal{N}} b_{Nn}(t_{i+1}) + B_N(t_{i+1}) \right) - \left( \sum_{n \in \mathcal{N}} b_{Nn}(t_i) + B_N(t_i) \right) \quad (20)$$

Equivalently, the node can minimize the total fee cost, which comes from two sources: from lost fees because of dropped transactions<sup>11</sup>, and from fees paid for rebalancing operations:

$$\begin{aligned} L(t_i, r(t_i)) &= \int_{\tau \in (t_i, t_{i+1})} (f(A_{LR}(\tau) - S_{LR}(\tau)) + f(A_{RL}(\tau) - S_{RL}(\tau))) d\tau \\ &\quad + \sum_{n \in \mathcal{N}} (F_{swap}^{in}(r_n^{in}(t_i)) + F_{swap}^{out}(r_n^{out}(t_i))) \end{aligned} \quad (21)$$

---

<sup>11</sup>Note that we assume the node knows not only about the transactions that reach it, but also about the transactions that are supposed to reach it but never do because of insufficient remote balances. This is not strictly true in practice, but the node can approximate it by observing the transactions during an interval in which the remote balances are both big enough so that no incoming transaction would fail and create an estimate based on this observation.

The two objectives at each timestep sum to  $\int_{\tau \in (t_i, t_{i+1})} (f(A_{LR}(\tau)) + f(A_{RL}(\tau))) d\tau$  (the fees that would be collected by the node if the total arriving amount had been processed), which is a quantity independent of the control action, and therefore maximizing the total fortune and minimizing the total fee cost are equivalent.

A control policy  $\pi = \{(t_i, r^\pi(t_i))\}_{i \in \mathbb{N}}$  consists of the times  $t_i$  and the corresponding actions  $r^\pi(t_i) = (r_L^{in}(t_i), r_L^{out}(t_i), r_R^{in}(t_i), r_R^{out}(t_i))$  taken from the set of allowed actions  $\mathcal{R} = [0, C_L]^2 \times [0, C_R]^2$ , and belongs to the set of admissible policies

$$\Pi = \left\{ \{(t_i, r(t_i))\}_{i \in \mathbb{N}} \text{ such that } r(t_i) \in \mathcal{R} \text{ for all } i \in \mathbb{N} \right\}$$

Ultimately, the goal of node  $N$  is to find a rebalancing policy that maximizes the long-term average expected fortune increase  $D$  (equivalently, minimizes the long-term average expected fee cost  $L$ ) over all admissible rebalancing policies:

$$\text{maximize}_{\pi \in \Pi} \lim_{H \rightarrow \infty} \frac{1}{t_H} \sum_{i=0}^H \mathbb{E} [D(t_i, r^\pi(t_i))] \quad (22)$$

subject to the constraints (8)–(12).

## 4 Heuristic and deep reinforcement learning-based rebalancing policies

In this section, we describe the steps we took in order to apply DRL to approximately solve the formulated MDP. We first outline two heuristic policies, which we will use later to benchmark our DRL-based solution.

### 4.1 Heuristic policies

---

**Algorithm 1:** Autoloop rebalancing policy

---

**Input:** *state* as in (7)  
**Parameters:**  $T_{check}$ , *low*, *high*

```

1 every  $T_{check}$  do
2   foreach neighbor  $n \in \mathcal{N}$  do
3      $midpoint = C_n \cdot (low + high) / 2$ 
4     if  $b_{Nn} < low \cdot C_n$  then
5       Swap-in amount =  $midpoint - b_{Nn}$ 
6     else if  $b_{Nn} > high \cdot C_n$  then
7       Swap-out amount =  $b_{Nn} - midpoint$ 
8     else
9       Perform no action

```

---

Autoloop [10, 11] is a policy that allows a node to schedule automatic swap-ins (resp. swap-outs) if its local balance falls below a minimum (resp. rises above a maximum) threshold expressed as a percentage of the channel’s capacity (Alg. 1)<sup>12</sup>. We would expect Autoloop to be suboptimal with respect to profit maximization in certain cases, as it does not take the expected demand into account and thus possibly performs rebalancing at times when it is not necessary.

This motivates us to define another heuristic policy that incorporates the empirical demand information. We call this policy Loopmax (Alg. 2), as its goal is to rebalance with the maximum possible amount and as infrequently as possible (without sacrificing transactions), based on the demand at each time. Loopmax keeps track of the total arriving amounts, and estimates the net change of each balance per unit time using the difference of the total amounts that arrived in each direction:

$$\hat{A}_{LN}^{net}(\tau) = -\hat{A}_{NL}^{net}(\tau) \triangleq \frac{1}{\tau} \int_{t \in [0, \tau]} \left( A_{RL}(t) - f(A_{RL}(t)) - A_{LR}(t) \right) dt \quad (23)$$

$$\hat{A}_{RN}^{net}(\tau) = -\hat{A}_{NR}^{net}(\tau) \triangleq \frac{1}{\tau} \int_{t \in [0, \tau]} \left( A_{LR}(t) - f(A_{LR}(t)) - A_{RL}(t) \right) dt \quad (24)$$

We first calculate the estimated time to depletion (*ETTD*) or saturation (*ETTS*) of the channel, depending on the direction of the net demand and the current balances, and using this time we dispatch a swap of the appropriate type not earlier than  $T_{check} + T_{conf}$  before depletion/saturation, and of the maximum possible amount. The rationale is that if e.g.  $ETTD \geq T_{check} + T_{conf}$ , the policy can leverage this fact to postpone starting a rebalancing until the next check time, since until then no transactions will have been dropped. If  $ETTD < T_{check} + T_{conf}$  though, the policy should act now, as otherwise it will end up dropping transactions during the following  $T_{check} + T_{conf}$  time. The maximum possible swap-out is constrained by the local balance at that time, while the maximum possible swap-in is constrained by the remote balance at that time<sup>13</sup> and the on-chain amount: an on-chain amount of  $B_N$  can support (by including fees) a net swap-in amount of at most  $\phi^{-1}(B_N)$ .

Compared to Autoloop, Loopmax has the advantage that it rebalances only when it is absolutely necessary and can thus achieve savings in swap fees. On the other hand, Loopmax’s aggressiveness can lead it to take extreme rebalancing decisions when the traffic is quite skewed in a particular direction (e.g. it can do a swap-in of almost the full capacity, which is very likely to fail due to randomness in the transaction arrivals). A small modification we can use on top of Alg. 2 to alleviate this is to define certain safety margins of liquidity that Loopmax should always leave intact on each side of the channel, so that incoming transactions do not find it depleted due to a large pending swap.

---

<sup>12</sup>The original Autoloop algorithm defines the thresholds in terms of inbound liquidity in a node’s channel. We adopt an equivalent balance-centric view instead.

<sup>13</sup>Actually, it is constrained by the remote balance at the time of the swap-in’s completion. We will improve this later using estimates of future balances.

---

**Algorithm 2:** Loopmax rebalancing policy

---

**Input:** *state* as in (7)  
**Parameters:**  $T_{check}$

```
1 every  $T_{check}$  do
2   Update  $\{\hat{A}_{Nn}^{net}\}_{n \in \mathcal{N}}$  according to (23)-(24)
3   foreach neighbor  $n \in \mathcal{N}$  do
4     if  $\hat{A}_{Nn}^{net} < 0$  then
5        $ETTD = b_{Nn} / |\hat{A}_{Nn}^{net}|$  /* estimated time to depletion */
6       if  $ETTD < T_{check} + T_{conf}$  then
7         Swap-in amount =  $\max\{\phi^{-1}(B_N), b_{nN}\}$  /* maximum possible swap in
8         */
9       else
10        Perform no action
11    else if  $\hat{A}_{Nn}^{net} > 0$  then
12       $ETTS = b_{nN} / \hat{A}_{Nn}^{net}$  /* estimated time to saturation */
13      if  $ETTS < T_{check} + T_{conf}$  then
14        Swap-out amount =  $b_{nN}$  /* maximum possible swap out */
15      else
16        Perform no action
17    else
18      Perform no action
```

---

## 4.2 Deep reinforcement learning algorithm design

Having formulated the problem as an MDP, we now need to find an (approximately) optimal policy. The problem is quite challenging for a number of reasons:

- The problem dynamics are not linear.
- The state and action spaces are continuous and thus tabular approaches are not applicable.
- There are time-dependent constraints on the actions.
- Choosing to not rebalance at a specific time requires special treatment, as otherwise the zero action will be sampled from a continuous action space with zero probability.

To tackle these challenges, we resort to approximate methods, and specifically Reinforcement Learning (RL). In the standard RL framework, an agent makes decisions based on a policy that is represented as a probability distribution over states and actions:  $p : p(s, a) \rightarrow [0, 1]$ , with  $p(s, a)$  being the probability that action  $a$  will be taken when the environment is in state  $s$ . Since our problem has continuous state and action spaces and the policy cannot be stored in tabular form, we need to use function approximation techniques. Neural networks

serve well the role of function approximators in many applications [12]. Some algorithms appropriate for this type of problems are Deep Deterministic Policy Gradient (DDPG) [13] and Soft Actor-Critic (SAC) [14]. We decided to use the latter because DDPG is known to exhibit extreme brittleness and hyperparameter sensitivity [15].

We now describe our methodology around how we engineer our DRL algorithm based on the vanilla SAC in order to arrive at a solution that deals with all the above challenges.

For the RL agent’s environment, we use as state the five balances (off- and on-chain) and the estimates of the remote balances at the time of the swap completion, each of them normalized appropriately: by the respective channel’s capacity, or by a total target fortune in the on-chain amount’s case. Thus, our state space is  $[0, 1]^7$ . As actions, instead of the 4-tuple of Section 3, we use a 2-tuple  $(r_L, r_R)$ , i.e. a single variable for each channel that can take both positive (swap-in) and negative (swap-out) values. Raw actions are sampled from the entire continuous action space; before the raw action is applied, it undergoes some processing described in the sequel.

As mentioned, an action with a coordinate equal to zero would be selected with zero probability. In reality, though, performing zero rebalancing in a channel when a swap is not necessary is important for maximizing the fortune/minimizing the costs, and something we would like the agent to learn to do. To this end, if the raw action is less than a threshold  $\rho_0$  (e.g. 20%) of the channel capacity, we force the applied action coordinate to be zero. This way, we make the zero action selectable with positive probability, and at the same time prevent the agent from performing swaps too small in size (which would increase the cost).

Moreover, in order to guide the algorithm to respect the constraints, we perform an additional processing step. The vanilla SAC algorithm [14] operates on an action space that is a compact subset of  $\mathbb{R}^k$  for all decision times. In our case, though, the allowed actions vary due to the time-dependent constraints (8)–(12). We therefore define the action space to be  $[-1, 1]^2$ , where each coordinate denotes the percentage not of the entire channel capacity, but only of the maximum amount available for the respective type of swap at that moment. We now focus on deriving these maximum amounts from the constraints.

All constraints are decoupled per channel, except for (12). However, we observe that given some traffic, mostly in the  $L$ -to- $R$  direction or mostly in the  $R$ -to- $L$  direction or equal in both directions, the local balances of node  $N$  will either deplete in one channel and accumulate in the other, or accumulate in both, but never both deplete. Thus, a swap-in in both channels in general will not be a good action. Therefore, for the RL solution’s purposes we can split (12) into two constraints, one for each channel, with the right-hand side of each being the entire amount  $B_N(t_i)$ . In case the agent does take the not advisable decision of swap-ins in both channels and their sum exceeds the on-chain amount, one of the two will simply fail.

Another useful observation is that when a swap-in is about to complete time  $T_{conf}$  after it was commenced, the remote balance in the respective channel needs to suffice (otherwise the swap-in will fail and a refund will be triggered as in eq. (19)):

$$r_n^{in}(t_i) \leq b_{nN}(t_i) + d_{nN}^{(t_i, t_i + T_{conf})} \quad \text{for all } i \in \mathbb{N}, n \in \mathcal{N} \quad (25)$$

Although (25) are not hard constraints when the decision is being made like the ones of Section 3.1.4, we would like to guide the agent to respect them. An obstacle is that the

swap-in decision is made at time  $t_i$ , when the node does not yet know the arriving amount  $d_{nN}^{(t_i, t_i + T_{conf})}$ . To approximate the right-hand side of (25) in terms of quantities known at time  $t_i$ , we can use the difference of the total (and not the successful as in  $d_{nN}$ 's definition) amounts that arrived in each direction (eq. (23)-(24)):

$$b_{nN}(t_i) + d_{nN}^{(t_i, t_i + T_{conf})} \approx \hat{b}_{nN}(t_i + T_{conf}) \triangleq (\min\{b_{nN}(t_i) + \hat{A}_{nN}^{net} \cdot T_{conf}, C_n\})^+ \quad (26)$$

A better estimate can be obtained by using the empirical amounts that succeeded in either direction:

$$\hat{S}_{LR}(\tau) \triangleq \frac{1}{\tau} \int_{t \in [0, \tau]} S_{LR}(t) dt \quad \text{and} \quad \hat{S}_{RL}(\tau) \triangleq \frac{1}{\tau} \int_{t \in [0, \tau]} S_{RL}(t) dt \quad (27)$$

Then the amount  $\hat{S}_{LR}$  (resp.  $\hat{S}_{RL}$ ) will be flowing in the  $L$ -to- $R$  (resp.  $R$ -to- $L$ ) direction for either the entire duration of  $T_{conf}$ , or until one of the balances in the respective direction is depleted:

$$\begin{aligned} \hat{b}_{LN}(t_i + T_{conf}) \triangleq & \left( \min \left\{ b_{LN}(t_i) - \hat{S}_{LR}(t_i) \min \left\{ T_{conf}, \frac{b_{LN}}{\hat{S}_{LR}(t_i)}, \frac{b_{NR}}{\hat{S}_{LR}(t_i)} \right\} \right. \right. \\ & \left. \left. + (1 - f_{prop}) \hat{S}_{RL}(t_i) \min \left\{ T_{conf}, \frac{b_{RN}}{\hat{S}_{RL}(t_i)}, \frac{b_{NL}}{\hat{S}_{RL}(t_i)} \right\}, C_L \right\} \right)^+ \end{aligned} \quad (28)$$

$$\begin{aligned} \hat{b}_{RN}(t_i + T_{conf}) \triangleq & \left( \min \left\{ b_{RN}(t_i) - \hat{S}_{RL}(t_i) \min \left\{ T_{conf}, \frac{b_{RN}}{\hat{S}_{RL}(t_i)}, \frac{b_{NL}}{\hat{S}_{RL}(t_i)} \right\} \right. \right. \\ & \left. \left. + (1 - f_{prop}) \hat{S}_{LR}(t_i) \min \left\{ T_{conf}, \frac{b_{LN}}{\hat{S}_{LR}(t_i)}, \frac{b_{NR}}{\hat{S}_{LR}(t_i)} \right\}, C_R \right\} \right)^+ \end{aligned} \quad (29)$$

Thus, the approximate version of (25) becomes:

$$r_n^{in}(t_i) \leq \hat{b}_{nN}(t_i + T_{conf}) \text{ for all } i \in \mathbb{N}, n \in \mathcal{N} \quad (30)$$

Note that we have given the agent more flexibility compared to Autoloop and Loopmax: it is allowed to perform swap-ins of amount bigger than the one allowed by the current balances, under the expectation that by the time of their completion the balances will be adequate.

Now we can write all constraints (8)–(12), (30) in terms of the 2-tuple  $(r_L, r_R)$  as follows:

$$r_n \in [-b_{Nn}, -\rho_{min}^{out}] \cup \left[ 0, \min\{\hat{b}_{nN}(t_i + T_{conf}), \phi^{-1}(B_N(t_i)), C_n\} \right], n \in \mathcal{N}$$

where  $\rho_{min}^{out} \triangleq M/(1 - F)$  is the minimum solution of (10).

If  $\rho_0 C_n \gg \rho_{min}^{out}$ , which should hold in practice as  $\rho_{min}^{out}$  is very small, we can write

$$r_n \in \left[ -b_{Nn}, \min\{\hat{b}_{nN}(t_i + T_{conf}), \phi^{-1}(B_N(t_i)), C_n\} \right], n \in \mathcal{N} \quad (31)$$

Table 1: Mapping of raw actions sampled from the learned distribution to final swap amounts requested for channel  $Nn$ ,  $n \in \mathcal{N}$

Raw action $r_n \in [-1, 1]$	Corresponding absolute amount $\tilde{r}_n$	Final requested swap amount
$r_n < 0$	$ r_n b_{Nn}$	swap out $\tilde{r}_n \mathbb{1}\{\tilde{r}_n \geq \rho_0 C_n\}$
$r_n \geq 0$	$r_n \min\{\hat{b}_{nN}(t_i + T_{conf}), \phi^{-1}(B_N(t_i)), C_n\}$	swap in $\tilde{r}_n \mathbb{1}\{\tilde{r}_n > \rho_0 C_n\}$

The final mapping of raw actions (sampled from the distribution on the entire action space) to the finally applied actions is shown in Table 1.

We craft the reward signal to guide the agent towards optimizing the objective: we add the node’s fortune increase (20) until the next check time, subtract the fee losses from transactions dropped until the next check time, and also subtract a fixed penalty for every swap the algorithm initiates and which eventually fails. A high-level sketch of the most important components of the final learning process described in this section is given in Alg. 3. We call the emerging policy “RebEL”: Rebalancing Enabled by Learning.

---

**Algorithm 3:** RL algorithm for RebEL policy

---

**Input:** *state* as in (7)  
**Parameters:**  $T_{check}$ , various learning parameters, penalty

- 1 **every**  $T_{check}$  **do**
- 2     Update estimates  $\hat{S}_{LR}$ ,  $\hat{S}_{RL}$  and  $\hat{b}_{LN}$ ,  $\hat{b}_{RN}$  according to (27)–(29)
- 3     Perform SAC gradient step to update policy distribution as in [14] based on replay memory
- 4     Fetch  $state \in [0, 1]^7$
- 5     Sample *rawAction* from  $[-1, 1]^2$  according to policy distribution
- 6     *processedAction* = *process*(*rawAction*) where *process*( $\cdot$ ) is described in Table 1
- 7     Apply *processedAction* and wait for its completion
- 8     *reward* = *fortuneAfter* – *fortuneBefore* – *lostFees* – *penalty* · *#OfFailedSwaps*
- 9     Fetch *nextState*  $\in [0, 1]^7$
- 10    Store transition (*state*, *rawAction*, *reward*, *nextState*) to replay memory

---

## 5 Evaluation

### 5.1 Simulator

In order to evaluate the performance of different rebalancing policies, we build a discrete event simulator of a relay node with two payment channels and rebalancing capabilities using Python SimPy [16]. The simulator treats each channel as a resource allowed to undergo at most one active swap at a time, and allows for parameterization of the initial balances, the transaction generation distributions (frequency, amount, number) for both sides of the channel, the different fees, the swap check and confirmation times, the rebalancing policy and the parameters of each policy<sup>14</sup>.

<sup>14</sup>The code is publicly available at <https://github.com/npapadis/payment-channel-rebalancing>.

## 5.2 Experimental setup

We simulate a relay node with two payment channels, each of a capacity of \$1000 split equally between the channel’s nodes. Transactions arrive from both sides as Poisson processes. We evaluate policies Autoloop, Loopmax and RebEL defined in Sec. 4.2, as well as the *None* policy that never performs any rebalancing. We use  $T_{check} = T_{conf} = 10$  minutes, miner fee  $M = \$2$ /on-chain transaction (tx), swap fee  $F = 0.5\%$ , 0.3 and 0.7 as the low and high liquidity thresholds of Autoloop, and 2 minutes worth of estimated traffic as safety margins for Loopmax. We run all experiments on a regular consumer laptop.

We experimented with different hyperparameters for the original SAC algorithm<sup>15</sup> as well as for RebEL parameters and reward shapes, and settled with the ones shown in Appendix B. We performed experiments for the transaction amount distribution being Uniform in  $[0, 50]$  and Gaussian with mean 25 and standard deviation 20, and the results were very similar. Therefore, all plots shown below are for the Gaussian amounts.

## 5.3 Results

### 5.3.1 The role of fees

Current median fee rates for transaction forwarding are in the order of  $3 \cdot 10^{-5}$  (\$/\$) or 0.003%<sup>16</sup>, while swap server fees are in the order of 0.5%<sup>17</sup> and miner fees are in the order of 2 \$/tx<sup>18</sup>. In order to see if a relay node can make a profit with such fees, we perform the following back-of-the-envelope calculation: A swap-in of amount  $r$  will cost the node  $rF + M$  in fees and will enable traffic of at most value  $r$  to be processed, which will yield profits  $r f_{prop}$  from relay fees. Therefore, the swap-in cannot be profitable if  $rF + M \geq r f_{prop}$ . Solving this inequality, we see that no positive amount  $r$  can be profitable if  $f_{prop} \leq F$ , while if  $f_{prop} > F$  a necessary (but not sufficient) condition for profitability is  $r > M/(f_{prop} - F)$ . The respective inequality for a swap-out of amount  $r$  is  $r - \frac{r-M}{1+F} \geq r f_{prop}$ , which shows that for  $f_{prop} \leq \frac{F}{1+F}$  no amount can be profitable and for  $f_{prop} > \frac{F}{1+F}$  a necessary condition for profitability is that  $r > \frac{M}{f_{prop}(1+F)-F}$ . With the current fees, we are in the non-profitable regime. Although the above inequalities are short-sighted in that they focus only on a specific action time, they do confirm the observation made by both the Lightning and the academic communities [17] that in order for relay nodes to be a profitable business, relay fees have to increase.

We now perform an experiment confirming this finding with the currently used fee values. We simulate a workload of demand in the  $L$ -to- $R$  direction: 60000  $L$ -to- $R$  and 15000  $R$ -to- $L$  transactions under a high (10 tx/minute  $L$ -to- $R$ , 2.5 tx/minute  $R$ -to- $L$ ) and a low (1 tx/minute  $L$ -to- $R$ , 0.25 tx/minute  $R$ -to- $L$ ) intensity. The node’s total fortune over time for high and low intensity are shown in Figs. 6(a) and 6(c) respectively. We see that regardless of the (non-*None*) rebalancing policy, the node’s fortune decreases over time, because rebalancing fees surpass any relay profits, which are small because of the small  $f_{prop}$  compared to  $F$ . In this regime, the node is better off not rebalancing at all. Still, our

<sup>15</sup>We used the PyTorch implementation in <https://github.com/pranz24/pytorch-soft-actor-critic>.

<sup>16</sup><https://1ml.com/statistics>

<sup>17</sup><https://lightning.engineering/loop>

<sup>18</sup>[https://ycharts.com/indicators/bitcoin\\_average\\_transaction\\_fee](https://ycharts.com/indicators/bitcoin_average_transaction_fee)

RebEL policy manages to learn this fact and after some point exhibits the desired behavior and stops rebalancing as well. Autoloop and Loopmax keep trying to rebalance and end up exhausting their entire on-chain balance, so the total fortune under them gets stuck after some point.

Taking a higher level view, we also conduct multiple experiments with the same demand as before but now while varying  $f_{prop}$ . The results of the total final fortune of each experiment (run for the same total time and averaged over 10 runs; error bars show the maximum and minimum values)

are shown in Fig. 6(b) under high demand and in Fig. 6(d) under low demand. We see that no rebalancing policy is profitable (i.e. better than *None*) as long as  $f_{prop} < 0.5\% = F$ , which confirms our back-of-the-envelope calculation. For higher values of  $f_{prop}$ , the node can make a profit. Although RebEL performs better for  $f_{prop} = 1\%$  for reasons discussed in Sec. 5.3.2, Autoloop and Loopmax sometimes perform better for even higher (and thus even farther from the current) fees, because the RebEL policy used

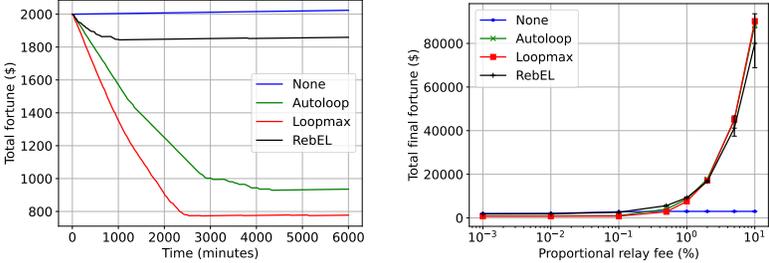
in this experiment is the one we tuned to operate best for the experiments of the next section that use  $f_{prop} = 1\%$ . In principle though, with different tuning, RebEL could outperform the other policies for higher values of  $f_{prop}$  as well.

in this experiment is the one we tuned to operate best for the experiments of the next section that use  $f_{prop} = 1\%$ . In principle though, with different tuning, RebEL could outperform the other policies for higher values of  $f_{prop}$  as well.

### 5.3.2 The role of the demand structure

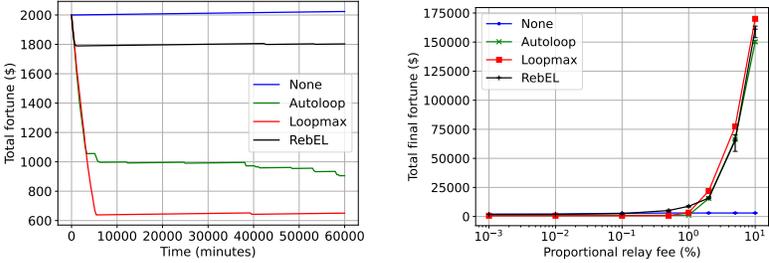
We now stay in the fee regime of possible profitability, i.e. by keeping  $f_{prop} = 1\%$ , and try to understand the role of the demand on the performance of the different policies. The results for the same high and low workload of skewed demand in the *L-to-R* direction as before are shown in Figs. 7 and 8.

RebEL outperforms all other policies under both demand regimes (Figs. 7(a), 8(a)), as it manages to strike a balance in terms of frequency and amount of rebalancing and transaction fee profits. This happens in a few 10-minute iterations under high demand



(a) Total fortune over time under high demand skewed in the *L-to-R* direction

(b) Total final fortune under high demand skewed in the *L-to-R* direction for different values of the proportional relay fee  $f_{prop}$



(c) Total fortune over time under low demand skewed in the *L-to-R* direction

(d) Total final fortune under low demand skewed in the *L-to-R* direction for different values of the proportional relay fee  $f_{prop}$

Figure 6: Experiments with different proportional relay fee  $f_{prop}$

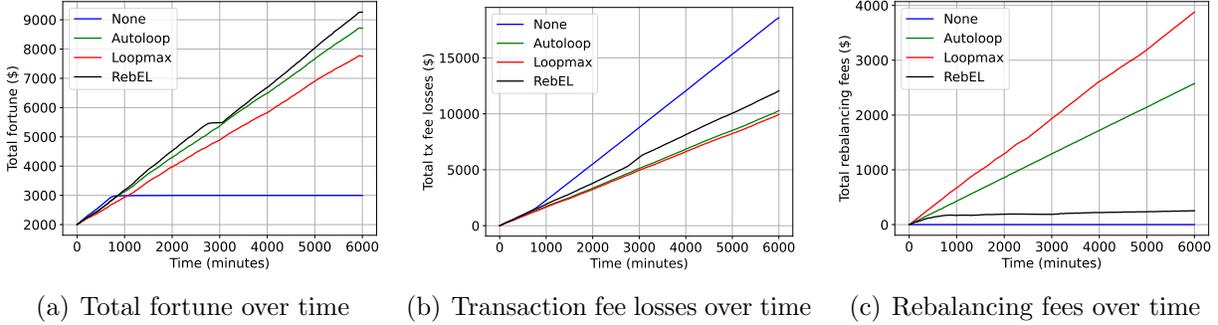


Figure 7: Total fortune, transaction fee losses and rebalancing fees over time under high demand skewed in the  $L$ -to- $R$  direction

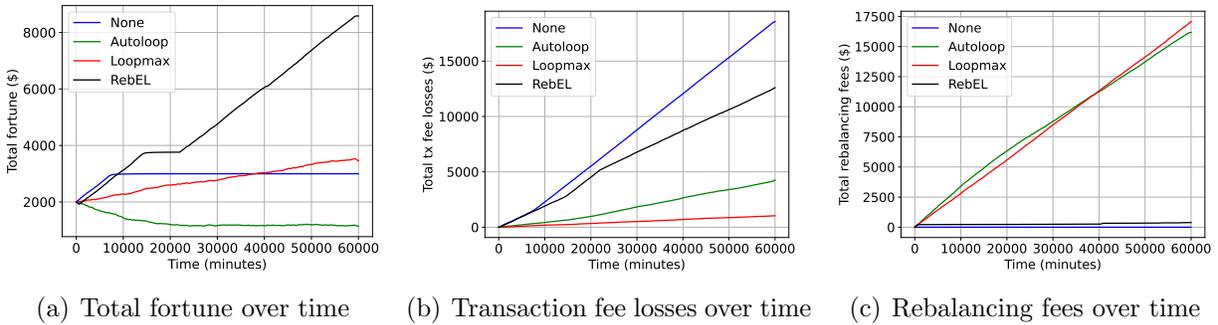


Figure 8: Total fortune, transaction fee losses and rebalancing fees over time under low demand skewed in the  $L$ -to- $R$  direction

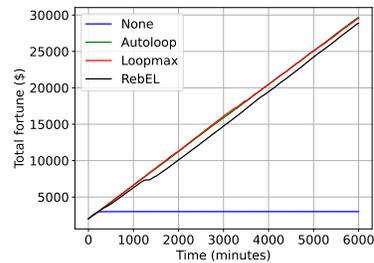
(corresponding to a few hours in real time), because balance changes are more pronounced in this case and help RebEL learn faster, while it takes about 1200 iterations under low demand, translating in 8.3 days of training. Both these training times should be reasonable for a relay node investing its capital to make a profit. We see that under both regimes the system without rebalancing (*None* policy) at some point reaches a state where almost all the balances are accumulated locally and no transactions can be processed anymore (hence the flattening in the *None* curve). Under high demand, Autoloop and Loopmax rebalance a lot (Fig. 7(c)) in order to minimize transaction fee losses (Fig. 7(b)), while RebEL sacrifices some transactions to achieve higher total fortune. Under low demand, RebEL rebalances only when necessary (Fig. 8(c)), even if this means sacrificing many more transactions (Fig. 8(b)), because simply rebalancing is not worth it at that low demand regime, in the sense that the potential profits during the 10-minute rebalancing check times are too low to justify the frequent rebalancing operations that the other policies apply. Loopmax eventually achieves a profit (although much lower than RebEL) because it tends to rebalance with higher amounts. On the contrary, Autoloop rebalances with small amounts, thus incurring significant costs from constant miner fees and eventually even making a loss compared to the initial node's fortune (Fig. 8(a)). Under high demand, there is a point around time 2700 where RebEL stalls for a bit, and the same happens under low demand between times 14000-22000. Upon more detailed inspection, this happens because all balances temporarily accumulate on the

local sides of the channels. RebEL takes some steps to again bring the channels to some balance (either actively by making a swap or passively by letting transactions flow) and subsequently completely recovers.

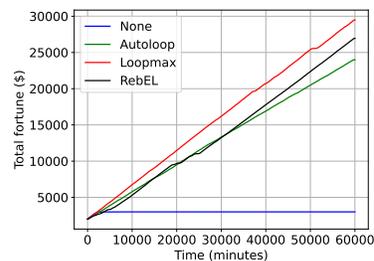
We now explore the special case of equal demands, by applying 60000 transactions arriving on each side in high (10 tx/minute) and low (1 tx/minute) intensity. Tuning some hyperparameters and making the penalty for failed swaps non-zero as shown in Appendix B gave better results for even demand specifically, so we use this configuration for the results of Fig. 9. We observe that all policies (except *None*) achieve higher total fortunes than before. This happens because the almost even traffic automatically rebalances the channel to some extent and therefore more fees can be collected in both directions and for larger amounts of time before the channels get stuck. RebEL is not as good for even traffic, because the net demand constantly oscillates around zero and this does not allow the agent to learn a good policy. It still manages though to surpass Autoloop pretty quickly under low demand, while if we run the simulation for longer times (not shown in the figure), we see that after time 78000 RebEL surpasses Loopmax as well. This translates to about 54 days of operation, which is a big time interval in practice, but is justified by the fact that the traffic is low and therefore more time is needed in order for the node to make a profit. However, even demand from both sides is a special case that is not likely to occur in practice, as usually the traffic follows some patterns, e.g. from clients to merchants. So the skewed demand scenario, where RebEL is superior, is also the most natural.

### 5.3.3 The role of initial conditions

We now examine how the initial conditions (capacities, balances) affect the performance. We evaluate all rebalancing policies for the skewed demand in the *L*-to-*R* direction scenario as before, but this time for channels of uneven capacities or initial balances. The results for high and low demand are shown in Figs. 10(a) and 11(a) respectively for  $C_L = 1000$ ,  $C_R = 500$  and the initial balances evenly distributed, in Figs. 10(b) and 11(b) respectively for  $C_L = 500$ ,  $C_R = 1000$  and the initial balances evenly distributed, and in Figs. 10(c) and 11(c) respectively for  $C_L = C_R = 1000$  but  $b_{NL} = b_{NR} = 1000$  (and so  $b_{LN} = b_{RN} = 0$ ). We see that RebEL performs well in all these cases as well. Depending on the exact arriving transactions, the little plateaus of RebEL happen at different points in time for the same reason as before, but in the end the learning algorithm recovers.

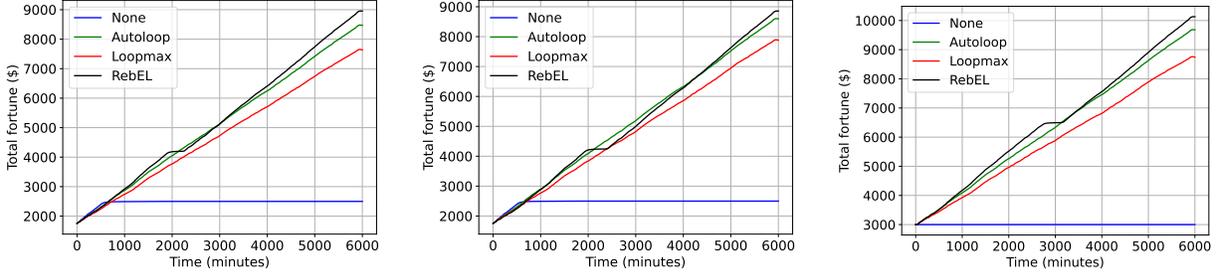


(a) High demand



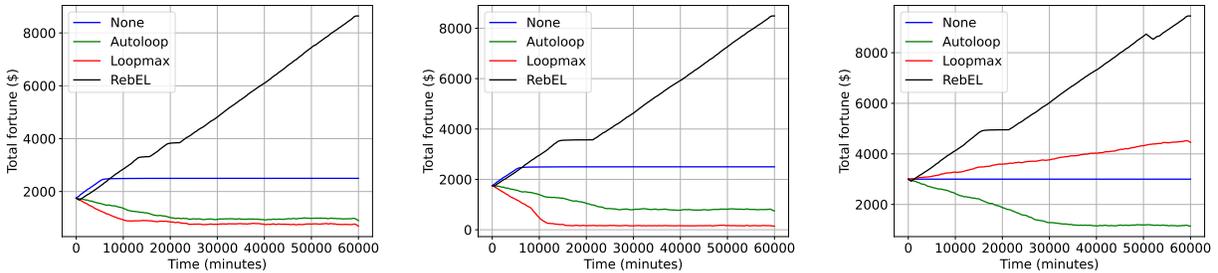
(b) Low demand

Figure 9: Total fortune over time under equal demand intensity from both sides



(a) Total fortune over time when  $C_L = 1000, C_R = 500$  (b) Total fortune over time when  $C_L = 500, C_R = 1000$  (c) Total fortune over time when initial balances are only local

Figure 10: Total fortune, transaction fee losses and rebalancing fees over time under high demand skewed in the  $L$ -to- $R$  direction for different initial conditions



(a) Total fortune over time when  $C_L = 1000, C_R = 500$  (b) Total fortune over time when  $C_L = 500, C_R = 1000$  (c) Total fortune over time when initial balances are only local

Figure 11: Total fortune, transaction fee losses and rebalancing fees over time under low demand skewed in the  $L$ -to- $R$  direction for different initial conditions

## 6 Discussion and future work

We now make some remarks on the design and the practical applicability of our DRL-based policy and discuss future extensions of our work.

*Design choices:* The objectives of Section 3.2 were defined as long-term expected average ones in order to match what a relay node would intuitively want to optimize, while the SAC algorithm works for long-term discounted objectives with a discount factor (usually set very close to 1), and including a maximum entropy term to enhance exploration<sup>19</sup>. We expect this difference to not be significant, and indeed the results show that the SAC-based policy performs well in practice. Furthermore, in Sec. 5 we presented results for specific parameters and rewards for the RL algorithm. Further tuning specific to the demand regime might lead to even higher returns for the RebEL policy. Additionally, improving the estimates of future balances by having the agent perform a “mini-simulation” of the transactions arriving in the following time interval based on past statistics could help the policy produce more informed decisions. Techniques from Model Predictive Control could also be applied [18].

Theoretically, a class of policies that could result in even higher fortune than the class

<sup>19</sup>The exact formula for the SAC objective can be found in Appendix A of [14].

(3.2) would be one that would allow rebalancing to happen at any point in continuous time instead of periodically. Optimization in such a model however would be extremely difficult, as an action taken now would affect the state both now and in the future (when rebalancing completes). Considering that practical policies like Autoloop applied today only check for rebalancing periodically, we follow the same path for the sake of tractability.

*Practical applicability:* An actual PCN node can apply our techniques as follows: the node may use our simulator with samples from its past demand, and try to tune the RL parameters and the reward to get better performance than the heuristic policies we defined or the one it is currently using; then, it will apply the policy learned in the simulator environment to the real node. Alternatively, a node may not use a simulator at all and directly learn a pre-parameterized policy on the fly from the empirical transaction data. In either case, the node can do occasional retraining with updated data to account for time-variance in the distribution of the arriving demand.

*Future directions:* Our two-channel DRL solution was a proof of concept that DRL can indeed be applied for profit management in PCNs. Armed with this knowledge, in future research we intend to study the more general case of a node being the center of a star graph of channels and trying to make a profit while rebalancing all of them appropriately. Another extension would be to allow the node to batch rebalancing operations into one on-chain transaction to save on on-chain fees. Moreover, in our work we considered the neighboring nodes  $L$  and  $R$  to be passive. Future work can investigate a game-theoretic framework where all PCN nodes are rational and compete against each other towards making a profit. Finally, it would also be interesting to compare the performance of different rebalancing methods, depending on the demand and channel conditions.

## 7 Related work

*Rebalancing methods:* Rebalancing via payments from a node to itself via a circular path of channels has been studied by [19–24]. Some of them take relay fees into account as we did, and some do not. [25] in particular performs circular rebalancing coupled with a rerouting scheme based on a metric that accounts for the average demand in a simple way (we did so too in defining the different balance estimates). [26,27] describe fee strategies that incentivize the balanced use of payment channels. [9] uses a game-theoretic lens to study the extent to which nodes can pay lower transaction fees by waiting patiently and reordering transactions instead of pursuing maximum efficiency. Perhaps the only work on submarine swaps, [28], shows that there is a possibility of liquidity arbitrage of Lightning liquidity providers by users, which in turn determines a market rate for acquiring liquidity, and then develops fee structures for properly pricing liquidity without overcharging regular users. In [29], a more holistic view is attempted regarding an optimization decision a blockchain node with an initial budget has to make: how to maximize the average gain per incoming transaction from a *known* distribution by choosing which channels to open, with what capacities and with what fees. However, the model ignores the channel opening costs by assuming it is possible to extend the channel’s lifetime arbitrarily, without though detailing how this would be done (e.g. via rebalancing). A recent development similar to submarine swaps is PeerSwap [30,31]: instead of buying funds from an LSP, a node can exchange funds on-/off-chain with its channel

neighbor directly. *Splicing* is another mechanism that replaces a channel with a new one with a different capacity while allowing transactions to flow in the meantime [32].

*Techniques:* Stochastic modeling and optimization in the blockchain space has been used both in layer-1 [33–36] for performance characterization, and in layer-2 for routing [37] and scheduling [38] of payments. Deep Reinforcement Learning has been broadly applied to approximately solve challenging optimization problems from various areas and to build systems that learn to manage resources directly from experience. For example, [12] applies DRL to the resource allocation problem of packing tasks under multiple resource demands, while [39] describes a DRL framework for solving a complex MDP underlying the incentives around selfish mining attacks in Bitcoin-like blockchains. Our profitable rebalancing problem resembles problems appearing in stochastic inventory control, without or with a positive lead time for replenishment. The so-called  $(s, S)$  threshold policies (if inventory level  $x < s$ , order  $S - x$ ; if  $x > s$ , do not order) can be proved to be optimal in certain settings [40–42]. Autoloop resembles these policies; however, our problem presents additional complexities due to the fact that there are more than one channels, with the balances of each affecting transaction processing in the other, leading us to a DRL-based approach. (Deep) RL has been applied extensively to inventory management problems as well [43, 44], although usually extensive tuning is necessary [45].

## 8 Conclusion

In this paper, we studied the problem of relay node profit maximization using submarine swaps, and demonstrated the feasibility of applying state-of-the-art DRL techniques for solving it, with our experiments showing that a SAC-based policy can outperform heuristic policies in most cases. We hope that this research will inspire further interest in designing capital management strategies in the complex world of PCNs based on learning from experience as an alternative to currently applied heuristics, and will be a step towards guaranteeing the profitability of the relay nodes and, consequently, the viability and scalability of the PCNs they sustain.

## References

- [1] K. Croman, C. Decker, I. Eyal, A. E. Gencer, A. Juels, A. Kosba, A. Miller, P. Saxena, E. Shi, E. Gün Sirer, D. Song, and R. Wattenhofer, “On scaling decentralized blockchains,” in *Financial Cryptography and Data Security*, J. Clark, S. Meiklejohn, P. Y. Ryan, D. Wallach, M. Brenner, and K. Rohloff, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2016, pp. 106–125.
- [2] J. Poon and T. Dryja, “The Bitcoin Lightning Network: scalable off-chain instant payments,” 2016. [Online]. Available: <https://lightning.network/lightning-network-paper.pdf>

- [3] L. Gudgeon, P. Moreno-Sanchez, S. Roos, P. McCorry, and A. Gervais, “SoK: Layer-two blockchain protocols,” in *Financial Cryptography and Data Security*, J. Bonneau and N. Heninger, Eds. Cham: Springer International Publishing, 2020, pp. 201–226.
- [4] N. Papadis and L. Tassiulas, “Blockchain-based payment channel networks: Challenges and recent advances,” *IEEE Access*, vol. 8, pp. 227 596–227 609, 2020. [Online]. Available: <https://doi.org/10.1109/ACCESS.2020.3046020>
- [5] A. Bosworth, “Submarine swaps on the Lightning Network,” 2018, [Online; accessed 9-August-2022]. [Online]. Available: <https://submarineswaps.github.io>
- [6] “Submarine Swap,” 2021. [Online]. Available: <https://wiki.ion.radar.tech/tech/research/submarine-swap>
- [7] J. Jager, “Loop Out in-depth,” 2019. [Online]. Available: <https://blog.lightning.engineering/technical/posts/2019/04/15/loop-out-in-depth.html>
- [8] “Zero Base Fee graph,” 2022. [Online]. Available: <https://lnrouter.app/graph/zero-base-fee>
- [9] Q. Bai, Y. Xu, and X. Wang, “Understanding the benefit of being patient in payment channel networks,” *IEEE Transactions on Network Science and Engineering*, vol. 9, no. 3, pp. 1895–1908, 2022.
- [10] Lightning Labs, “Autoloop,” 2022. [Online]. Available: <https://github.com/lightninglabs/loop/blob/master/docs/autoloop.md>
- [11] Carla Kirk-Cohen, “Autoloop: Lightning Liquidity You Can Set and Forget!” 2020. [Online]. Available: <https://lightning.engineering/posts/2020-11-24-autoloop>
- [12] H. Mao, M. Alizadeh, I. Menache, and S. Kandula, “Resource management with deep reinforcement learning,” in *Proceedings of the 15th ACM Workshop on Hot Topics in Networks*, ser. HotNets ’16. New York, NY, USA: Association for Computing Machinery, 2016, p. 50–56. [Online]. Available: <https://doi.org/10.1145/3005745.3005750>
- [13] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, “Continuous control with deep reinforcement learning,” in *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*, Y. Bengio and Y. LeCun, Eds., 2016. [Online]. Available: <http://arxiv.org/abs/1509.02971>
- [14] T. Haarnoja, A. Zhou, K. Hartikainen, G. Tucker, S. Ha, J. Tan, V. Kumar, H. Zhu, A. Gupta, P. Abbeel, and S. Levine, “Soft actor-critic algorithms and applications,” *CoRR*, vol. abs/1812.05905, 2018. [Online]. Available: <http://arxiv.org/abs/1812.05905>
- [15] Y. Duan, X. Chen, R. Houthoofd, J. Schulman, and P. Abbeel, “Benchmarking deep reinforcement learning for continuous control,” in *Proceedings of the 33rd International Conference on International Conference on Machine Learning - Volume 48*, ser. ICML’16. JMLR.org, 2016, p. 1329–1338.

- [16] O. Lünsdorf and S. Scherfke, “SimPy,” 2022. [Online]. Available: <https://simpy.readthedocs.io>
- [17] F. Béres, I. A. Seres, and A. A. Benczúr, “A cryptoeconomic traffic analysis of Bitcoin’s Lightning Network,” *Cryptoeconomic Systems*, 6 2020. [Online]. Available: <https://cryptoeconomicssystemspubpub.org/pub/b8rb0ywn>
- [18] U. Rosolia and F. Borrelli, “Learning model predictive control for iterative tasks. A data-driven control framework,” *IEEE Transactions on Automatic Control*, vol. 63, no. 7, pp. 1883–1896, 2018.
- [19] R. Khalil and A. Gervais, “Revive: Rebalancing off-blockchain payment networks,” in *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security, CCS 2017, Dallas, TX, USA, October 30 - November 03, 2017*, B. M. Thuraisingham, D. Evans, T. Malkin, and D. Xu, Eds. ACM, 2017, pp. 439–453. [Online]. Available: <https://doi.org/10.1145/3133956.3134033>
- [20] R. Pickhardt and M. Nowostawski, “Imbalance measure and proactive channel rebalancing algorithm for the lightning network,” in *IEEE International Conference on Blockchain and Cryptocurrency, ICBC 2020, Toronto, ON, Canada, May 2-6, 2020*. IEEE, 2020, pp. 1–5. [Online]. Available: <https://doi.org/10.1109/ICBC48266.2020.9169456>
- [21] Z. Avarikioti, K. Pietrzak, I. Salem, S. Schmid, S. Tiwari, and M. Yeo, “HIDE & SEEK: Privacy-preserving rebalancing on payment channel networks,” *CoRR*, vol. abs/2110.08848, 2021. [Online]. Available: <https://arxiv.org/abs/2110.08848>
- [22] M. Xu, Y. Zhang, F. Xu, and S. Zhong, “Privacy-preserving optimal recovering for the nearly exhausted payment channels,” in *2021 IEEE/ACM 29th International Symposium on Quality of Service (IWQOS)*, 2021, pp. 1–10.
- [23] Z. Hong, S. Guo, R. Zhang, P. Li, Y. Zhan, and W. Chen, “Cycle: Sustainable off-chain payment channel network with asynchronous rebalancing,” in *2022 52nd Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DSN)*, 2022, pp. 41–53.
- [24] M. Bastankhah, K. Chatterjee, M. A. Maddah-Ali, S. Schmid, J. Svoboda, and M. Yeo, “Online admission control and rebalancing in payment channel networks,” 2022. [Online]. Available: <https://arxiv.org/abs/2209.11936>
- [25] N. Awathare, Suraj, Akash, V. J. Ribeiro, and U. Bellur, “Rebal: Channel balancing for payment channel networks,” in *2021 29th International Symposium on Modeling, Analysis, and Simulation of Computer and Telecommunication Systems (MASCOTS)*, 2021, pp. 1–8.
- [26] G. D. Stasi, S. Avallone, R. Canonico, and G. Ventre, “Routing payments on the lightning network,” in *IEEE International Conference on Internet of Things (iThings) and IEEE Green Computing and Communications (GreenCom) and*

- IEEE Cyber, Physical and Social Computing (CPSCoM) and IEEE Smart Data (SmartData), iThings/GreenCom/CPSCoM/SmartData 2018, Halifax, NS, Canada, July 30 - August 3, 2018.* IEEE, 2018, pp. 1161–1170. [Online]. Available: <https://doi.org/10.1109/Cybermatics.2018.2018.00209>
- [27] Y. van Engelshoven and S. Roos, “The merchant: Avoiding payment channel depletion through incentives,” in *IEEE International Conference on Decentralized Applications and Infrastructures, DAPPS 2021, Online Event, August 23-26, 2021.* IEEE, 2021, pp. 59–68. [Online]. Available: <https://doi.org/10.1109/DAPPS52256.2021.00012>
- [28] J. I. R. Echenique and N. Burtsey, “Pricing liquidity for Lightning wallets,” 2022, [Online; accessed 9-August-2022]. [Online]. Available: <https://github.com/GaloyMoney/liquidity-fees-paper>
- [29] O. Ersoy, S. Roos, and Z. Erkin, “How to profit from payments channels,” in *Financial Cryptography and Data Security - 24th International Conference, FC 2020, Kota Kinabalu, Malaysia, February 10-14, 2020 Revised Selected Papers*, ser. Lecture Notes in Computer Science, J. Bonneau and N. Heninger, Eds., vol. 12059. Springer, 2020, pp. 284–303. [Online]. Available: [https://doi.org/10.1007/978-3-030-51280-4\\_16](https://doi.org/10.1007/978-3-030-51280-4_16)
- [30] “PeerSwap,” 2022. [Online]. Available: <https://www.peerswap.dev>
- [31] W. Togami and K. Nick, “PeerSwap: Decentralized P2P LN Balancing Protocol,” 2021. [Online]. Available: <https://blockstream.com/assets/downloads/2021-11-16-PeerSwap-Announcement.pdf>
- [32] R. Russell, “Splicing Proposal,” 2018. [Online]. Available: <https://lists.linuxfoundation.org/pipermail/lightning-dev/2018-October/001434.html>
- [33] A. Dembo, S. Kannan, E. N. Tas, D. Tse, P. Viswanath, X. Wang, and O. Zeitouni, “Everything is a race and Nakamoto always wins,” in *Proceedings of the 2020 ACM SIGSAC Conference on Computer and Communications Security*, ser. CCS ’20. New York, NY, USA: Association for Computing Machinery, 2020, p. 859–878. [Online]. Available: <https://doi.org/10.1145/3372297.3417290>
- [34] P. Gaži, A. Kiayias, and A. Russell, “Tight consistency bounds for bitcoin,” in *Proceedings of the 2020 ACM SIGSAC Conference on Computer and Communications Security*, ser. CCS ’20. New York, NY, USA: Association for Computing Machinery, 2020, p. 819–838. [Online]. Available: <https://doi.org/10.1145/3372297.3423365>
- [35] N. Papadis, S. Borst, A. Walid, M. Grissa, and L. Tassiulas, “Stochastic models and wide-area network measurements for blockchain design and analysis,” in *IEEE INFOCOM 2018 - IEEE Conference on Computer Communications.* IEEE, apr 2018, pp. 2546–2554. [Online]. Available: <https://ieeexplore.ieee.org/document/8485982/>
- [36] J. Mišić, V. B. Mišić, X. Chang, S. G. Motlagh, and M. Z. Ali, “Modeling of Bitcoin’s blockchain delivery network,” *IEEE Transactions on Network Science and Engineering*, vol. 7, no. 3, pp. 1368–1381, 2020.

- [37] S. M. Varma and S. T. Maguluri, “Throughput optimal routing in blockchain-based payment systems,” *IEEE Transactions on Control of Network Systems*, vol. 8, no. 4, pp. 1859–1868, 2021.
- [38] N. Papadis and L. Tassiulas, “Payment Channel Networks: Single-Hop Scheduling for Throughput Maximization,” in *IEEE INFOCOM 2022 - IEEE Conference on Computer Communications*, 2022, pp. 900–909.
- [39] R. Bar-Zur, A. Abu-Hanna, I. Eyal, and A. Tamar, “WeRLman: To tackle whale (transactions), go deep (RL),” in *Proceedings of the 15th ACM International Conference on Systems and Storage*, ser. SYSTOR ’22. New York, NY, USA: Association for Computing Machinery, 2022, p. 148. [Online]. Available: <https://doi.org/10.1145/3534056.3535005>
- [40] H. E. Scarf, “The optimality of (S, s) policies in the dynamic inventory problem,” 1959.
- [41] Y.-S. Zheng, “A simple proof for optimality of (s, S) policies in infinite-horizon inventory systems,” *Journal of Applied Probability*, vol. 28, no. 4, pp. 802–810, 1991. [Online]. Available: <http://www.jstor.org/stable/3214683>
- [42] S. P. Sethi and F. Cheng, “Optimality of (s, S) policies in inventory models with markovian demand,” *Operations Research*, vol. 45, no. 6, pp. 931–939, December 1997. [Online]. Available: <https://ideas.repec.org/a/inm/oropre/v45y1997i6p931-939.html>
- [43] B. Van Roy, D. Bertsekas, Y. Lee, and J. Tsitsiklis, “A neuro-dynamic programming approach to retailer inventory management,” in *Proceedings of the 36th IEEE Conference on Decision and Control*, vol. 4, 1997, pp. 4052–4057 vol.4.
- [44] R. N. Boute, J. Gijsbrechts, W. van Jaarsveld, and N. Vanvuchelen, “Deep reinforcement learning for inventory control: A roadmap,” *Eur. J. Oper. Res.*, vol. 298, no. 2, pp. 401–412, 2022. [Online]. Available: <https://doi.org/10.1016/j.ejor.2021.07.016>
- [45] J. Gijsbrechts, R. N. Boute, J. A. V. Mieghem, and D. J. Zhang, “Can deep reinforcement learning improve inventory management? Performance on lost sales, dual-sourcing, and multi-echelon problems,” *Manuf. Serv. Oper. Manag.*, vol. 24, no. 3, pp. 1349–1368, 2022. [Online]. Available: <https://doi.org/10.1287/msom.2021.1064>

## A Example of channel depletion under symmetric demand

Symmetric demand on two endpoints of a multihop path can cause imbalance due to fees withheld by intermediate nodes. Fig. 12 shows the evolution over time of a subnetwork of three channels with symmetric demand of amount 20 arriving alternately from either side of the path. When each transaction is relayed by node  $B$ , a 50% fee is withheld and the remaining amount of 10 is forwarded to the next channel in the path. We see that even though the end-to-end path demand is symmetric, after a few steps the channels get unbalanced and stop being able to process any more transactions<sup>20</sup>.

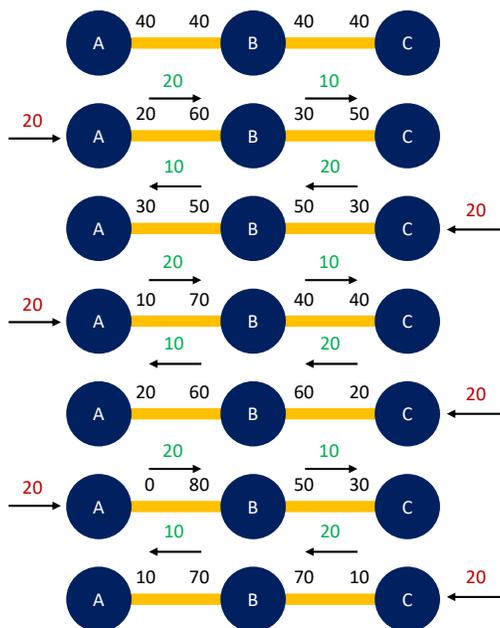


Figure 12: An example of a PCN getting stuck even though the demand is symmetric. Demand is shown in red, forwarded amounts after a 50% fee withholding are shown in green, and channel balances are shown in black.

<sup>20</sup>The 50% fee is not realistic and is only used for the purposes of this example. With the real much lower fees the channels will similarly get stuck after a larger number of steps.

## B Hyperparameters and rewards

Table 2: SAC hyperparameters used for the different experiments of Sec. 5

SAC hyperparameter	Parameter value for skewed demand experiments	Parameter value for even demand experiments
policy	Gaussian	
optimizer	Adam	
learning rate	0.0003	0.006
discount	0.99	
replay buffer size	$10^5$	
number of hidden layers (all neural networks)	2	
number of hidden units per layer	256	
number of samples per minibatch	10	
temperature	0.05	0.005
nonlinearity	ReLU	
target smoothing coefficient	0.005	
target update interval	1	
gradient steps	1	
automatic entropy tuning	False	True
initial random steps	10	

Table 3: Parameters used in ReBEL’s representation or processing of the states, actions, and rewards

ReBEL parameter	Parameter value for skewed demand experiments	Parameter value for even demand experiments
on-chain amount normalization constant	60	
minimum swap threshold $\rho_0$	0.2	
penalty per swap failure	0	10