FUTUREFILL: FAST GENERATION FROM CONVOLU TIONAL SEQUENCE MODELS

Anonymous authors

004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

023

Paper under double-blind review

ABSTRACT

We address the challenge of efficient auto-regressive generation in sequence prediction models by introducing FutureFill—a method for fast generation that applies to any sequence prediction algorithm based on convolutional operators. Our approach reduces the generation time requirement from quadratic to quasilinear relative to the context length. Additionally, FutureFill requires a prefill cache sized only by the number of tokens generated, which is smaller than the cache requirements for standard convolutional and attention-based models. We validate our theoretical findings with experimental evidence demonstrating correctness and efficiency gains in a synthetic generation task.

1 INTRODUCTION

Large Transformer models Vaswani et al. (2017) have become the method of choice for sequence prediction tasks such as language modeling and machine translation. Despite their success, they face a key computational limitation: the attention mechanism, their core innovation, incurs a quadratic computational cost during training and inference. This inefficiency has spurred interest in alternative architectures that can handle long sequences more efficiently.

Convolution-based sequence prediction models Li et al. (2022); Poli et al. (2023); Agarwal et al. (2023); Fu et al. (2024) have emerged as strong contenders, primarily due to their ability to leverage 031 the Fast Fourier Transform (FFT) for near-linear scaling with sequence length during training. These models build upon the advancements in State Space Models (SSMs), which have shown promise 033 in modeling long sequences across diverse modalities Gu et al. (2021a); Dao et al. (2022); Gupta 034 et al. (2022); Orvieto et al. (2023); Poli et al. (2023); Gu & Dao (2023). Convolutional models offer a more general framework than SSMs because they can represent any linear dynamical system (LDS) without being constrained by the dimensionality of hidden states Agarwal et al. (2023). This 037 flexibility has led to recent developments that theoretically and empirically handle longer contexts 038 more effectively. Notable among these are Spectral State Space Models or Spectral Transform Units (STUs) Agarwal et al. (2023), which use spectral filtering algorithms Hazan et al. (2017; 2018) to transform inputs into better-conditioned bases for long-term memory. Another approach is Hyena 040 Poli et al. (2023), which learns implicitly parameterized Markov operators. Both methods exploit 041 the duality between time-domain convolution and frequency-domain multiplication to accelerate 042 prediction via the FFT. 043

While SSMs and recurrent models benefit from fast inference times independent of sequence length,
making them attractive for large-scale language modeling, convolutional models have been hindered
by slower token generation during inference. The best-known result for generating tokens with
convolutional models is quadratic in sequence length—comparable to attention-based models (see
Massaroli et al. (2024) Lemma 2.1). This limitation has prompted research into distilling state-space
models from convolutional models Massaroli et al. (2024), but such approximations lack comprehensive understanding regarding their approximation gaps due to the broader representational capacity of convolutional models.

In this paper, we address the problem of exact auto-regressive generation from given convolutional
 models, significantly improving both the generation time and cache size requirements. We present
 our main results in two settings:

- 1. Generation from Scratch: When generating L tokens from scratch, we demonstrate that long convolutional sequence predictors can generate these tokens in total time $O(L \log^2 L)$ with total memory O(L). This improves upon previous methods that require $O(L^2)$ time for generation. We further provide a memory-efficient version wherein the total runtime increases to $O(L^{3/2}\sqrt{\log(L)})$ but the memory requirement is bounded by $O(\sqrt{L\log L})$.
 - 2. Generation with Prompt: When generating K tokens starting from a prompt of length L, we show that the total generation time is $O(L \log L + K \log^2 K)$ with a cache size requirement of O(K). Previously, the best-known requirements for convolutional models were a total generation time bounded by $O(L \log L + LK + K^2)$ and a cache size bounded by O(L) (Massaroli et al., 2024).

Importantly, our results pertain to provably exact generation from convolutional models without relying on any approximations. Moreover, our methods are applicable to any convolutional model, regardless of how it was trained. The following table compares our algorithm with a standard exact implementation of convolution. We also provide a comparison of the time and cache size requirements for exact computation in attention-based models. 069

| Method | Runtime | Memory |
|---------------------|------------------------|------------------|
| Standard Conv | L^2 | L |
| Standard Attn. | L^2 | L |
| EpochedFF (ours) | $L^{3/2}\sqrt{\log L}$ | $\sqrt{L\log L}$ |
| ContinuousFF (ours) | $L\log^2 L$ | L |

| Prefill+Genertation Runtime | Generation Cache Size |
|------------------------------------|--------------------------|
| $LK + L\log L + K^2$ | L |
| $L^2 + KL$ | L |
| $L \log L + K^{3/2} \sqrt{\log K}$ | K |
| $L\log L + K\log^2 K$ | K |

(a) Comparison for generating L tokens from scratch. Runtime is in asymptotic notation, i.e. $O(\cdot)$ is omitted for brevity.

(b) Comparison for generating K tokens starting from a prompt of length L, runtime and cachesize are in asymptotic notation, i.e. $O(\cdot)$ is omitted for brevity.

Our results for generation from convolutional models are based on building efficient algorithms for an online version of the problem of computing convolutions. In this problem, the algorithm is tasked to compute the convolution of two sequences $u * \phi$, however the challenge is to release iteratively at time t the value of $[u * \phi]_t$, where the sequence ϕ is fully available to the algorithm but the sequence u streams in one-coordinate at a time.

While the FFT algorithm allows for an $O(L \log L)$ -time offline algorithm for the convolution of two 088 L-length sequences, whether a similar result exists for the online model was not known. Naively, 089 since $[u*\phi]_t = \langle u_{1:t}, \phi_{t:1} \rangle$, the total output can be computed in time $O(L^2)$. In this paper we demon-090 strate using repeated calls to appropriately constructed FFT-subroutines to compute the *future* effect 091 of past tokens (a routine we call FutureFill), one can compute the convolution in the online model 092 with a total computational complexity of $O(L \log^2(L))$, nearly matching its offline counterpart and 093 significantly improving over the naive algorithm which was the best known (Massaroli et al., 2024). 094

It is worth noting that the naive algorithm for computing online convolution, albeit slow, does not 095 require any additional memory other than the memory used for storing the sequences v, w. Such 096 memory is often a bottleneck in practical sequence generation settings and is referred to as the size of the generation cache. For context the size of the generation cache for attention models is O(L), 098 i.e. proportional to the length of the prefill-context and the generation length. We further show that when generating from convolutional models, one can construct a trade-off for the computational 100 complexity (i.e. flops) and memory (i.e. generation cache size) using the FutureFill sub-routine. We 101 highlight two points on this trade-off spectrum via two algorithmic setups both employing FutureFill. 102 We detail this trade-off in Table 1.

103 104

105

055

056

058

060

061

062

063

064

065

066

067

068

071

073

079

080

081 082

083

084

085

087

1.1 RELATED WORK

State space models and convolutional sequence prediction. Recurrent neural networks have 106 been revisited in the recent deep learning literature for sequential prediction in the form of state 107 space models (SSM), many of whom can be parameterized as convolutional models. Gu et al. (2020)

| 108 | Algorithm | Computational | Memory (Generation |
|-----|------------------------------|--------------------|--------------------|
| 109 | | Complexity (Flops) | Cache-Size) |
| 110 | Naive | $O(L^2)$ | O(1) |
| 111 | Epoched-FutureFill (ours) | $O(L^{3/2}\log L)$ | $O(\sqrt{L})$ |
| 112 | Continuous-FutureFill (ours) | $O(L \log^2 L)$ | O(L) |

Table 1: Comparison of results for Online convolution.

propose the HiPPO framework for continuous-time memorization, and shows that with a special 117 class of system matrices A (HiPPO matrices), SSMs have the capacity for long-range memory. Later 118 work Gu et al. (2021b;a); Gupta et al. (2022); Smith et al. (2023) focus on removing nonlinearities 119 and devising computationally efficient methods that are also numerically stable. To improve the 120 performance of SSMs on language modeling tasks Dao et al. (2022) propose architectural changes 121 as well as faster FFT algorithms with better hardware utilization, to close the speed gap between 122 SSMs and Transformers. Further investigation in Orvieto et al. (2023) shows that training SSM is 123 brittle in terms of various hyperparameters. Various convolutional models have been proposed for 124 sequence modelling, see e.g. Fu et al. (2023); Li et al. (2022); Shi et al. (2023a). These papers 125 parameterize the convolution kernels with specific structures. The Hyena architecture was proposed in Poli et al. (2023) and distilling it into a SSM was studied in Massaroli et al. (2024). Other 126 studies in convolutional models include LongConv Fu et al. (2023) and SGConv Li et al. (2022) 127 architectures, as well as multi-resolution convolutional models Shi et al. (2023b). 128

Spectral filtering. A promising technique for learning in linear dynamical systems with long 130 memory is called spectral filtering put forth in Hazan et al. (2017). This work studies online pre-131 diction of the sequence of observations y_t , and the goal is to predict as well as the best symmetric 132 LDS using past inputs and observations. Directly learning the dynamics is a non-convex optimiza-133 tion problem, and spectral filtering is developed as an improper learning technique with an efficient, 134 polynomial-time algorithm and near-optimal regret guarantees. Different from regression-based 135 methods that aim to identify the system dynamics, spectral filtering's guarantee does not depend on 136 the stability of the underlying system, and is the first method to obtain condition number-free regret 137 guarantees for the MIMO setting. Extension to asymmetric dynamical systems was further studied 138 in Hazan et al. (2018). Spectral filtering is particularly relevant to this study since it is a convolutional model with fixed filters. Thus, our results immidiately apply to this technique and imply 139 provable regret bounds with guaranteed running time bounds in the online learning model which 140 improve upon state of the art. 141

142

129

113

114 115 116

Online learning and regret minimization in sequence prediction. The methodology of online 143 convex optimization, see e.g. Hazan et al. (2016), applies to sequences prediction naturally. In 144 this setting, a learner iteratively predicts, and suffers a loss according to an adversarially chosen 145 loss function. Since nature is assumed to be adversarial, statistical guarantees are not applicable, 146 and performance is measured in terms of regret, or the difference between the total loss and that 147 of the best algorithm in hindsight from a class of predictors. This is a particulary useful setting 148 for sequential prediction since no assumption about the sequence is made, and it leads to robust 149 methods. Sequential prediction methods that apply to dynamical systems are more complex as they 150 incorporate the notion of a state. Recently the theory of online convex optimization has been applied 151 to learning in dynamical systems, and in this context, the spectral filtering methodology was devised. See Hazan & Singh (2022) for an introduction to this area. 152

153 154

155

2 SETTING

156 2.1 Online Convolutions 157

158 **Notation:** For an input sequence $\{u_t\}$ we denote by $u_{1:t}$ the sequence of inputs u_1, \dots, u_t . For any 159 $i \leq j$ let $u_{i:j}$ denote the sub-sequence $u_i, u_{i+1}, \ldots u_j$. When $i > j, u_{i:j}$ denotes the subsequence $u_{i:i}$ in reverse order. Thus $u_{t:1}$ represents the sequence in reverse order. We also denote [k] =160 $\{1, 2, ..., k\}$ as a set of k natural numbers. Given a multi-dimensional sequence $u_1 \dots u_t$ where 161 each $u_i \in \mathbb{R}^d$ and given a vector $v \in \mathbb{R}^t$, for brevity of notation we overload the definition of inner products by defining $y = \langle v, u_{1:t} \rangle$ with $y \in \mathbb{R}^d$ as $y_j = \sum_{i=1}^t v_i \cdot [u_i]_j \in \mathbb{R}$. That is the inner-product along the time dimension is applied on each input dimension separately.

Convolution: The convolution operator between two vectors $u, \phi \in \mathbb{R}^t$ outputs a sequence of length t whose element at any position $s \in [t]^1$ is defined as

$$[u * \phi](s) = \sum_{i=1}^{s} u_i \phi_{s+1-i} = \langle u_{1:s}, \phi_{s:1} \rangle.$$
(1)

171 A classical result in the theory of algorithms is that given two vectors $u, \phi \in \mathbb{R}^t$, their convolution 172 can be computed in time $O(t \log t)$, using the FFT algorithm.

173 **Online Convolution:** We consider the problem of performing the convolution $u * \phi$ when one of 174 the sequences ϕ is fully available to the model, however the other sequence u streams in, i.e. the 175 element u_t is made available to the model at the start of round t, at which point it is expected to 176 release the output $[u * \phi]_t$. This model of online convolution is immediately relevant to the online 177 auto-regressive generation of tokens from a convolutional sequence model as the output token at time 178 t becomes the input for the next round and hence is only available post generation. In this setting, 179 the sequence u corresponds to generated tokens and the sequence ϕ corresponds to the convolutional 180 filter which the model has full access to. We further detail the setup of sequence generation in the 181 next subsection.

182 183

184

199 200 201

202

208

209

210 211

215

165

166

2.2 SEQUENCE PREDICTION:

In sequence prediction, the input is a sequence of tokens denoted $u_1, ..., u_t, ...,$ where $u_t \in \mathbb{R}^{d_{in}}$. The predictor's task is to generate a sequence $\hat{y}_1, ..., \hat{y}_t, ...,$ where $\hat{y}_t \in \mathbb{R}^{d_{out}}$ is generated after observing $u_1, ..., u_{t-1}$. The output y_t is observed after the predictor generates \hat{y}_t . The quality of the prediction is measured by the distance between the predicted and observed outputs according to a loss function $\ell_t(\hat{y}_t, y_t)$, for example the mean square error $\|\hat{y}_t - y_t\|^2$.

1902.3 ONLINE SEQUENCE PREDICTION

In the **online sequence prediction** setting, an online learner iteratively sees an input u_t and has to predict output \hat{y}_t , after which the true output y_t is revealed. The goal is to minimize error according to a given Lipschitz loss function $\ell_t(y_t, \hat{y}_t)$. In online learning it is uncommon to assume that the true sequence was generated by the same family of models as those learned by the learner. For this reason the metric of performance is taken to be regret. Given a class of possible predictors, the goal is to minimize regret w.r.t. these predictors. For example, a linear predictor predicts according to the rule

$$\pi_{M_{1:k},N_{1:l}}(u_{1:t},y_{1:t-1}) = \sum_{i=1}^{k} M_i u_{t-i} + \sum_{j=1}^{l} N_j y_{t-j}.$$

The performance of a prediction algorithm \mathcal{A} is measured by regret, or difference in total loss, vs. a class of predictors \prod , such as that of linear predictors, e.g.

$$\operatorname{Regret}_{T}(\mathcal{A}) = \sum_{t=1}^{T} \ell_{t}(y_{t}, \hat{y}_{t}^{\mathcal{A}}) - \min_{\pi \in \prod} \sum_{t=1}^{T} \ell_{t}(y_{t}, \hat{y}_{t}^{\pi}).$$

This formulation is valid for online sequence prediction of any signal. We are particularly interested in signals that are generated by dynamical systems. A time-invariant linear dynamical system is given by the dynamics equations

$$x_{t+1} = Ax_t + Bu_t + w_t, \ y_t = Cx_t + Du_t + \zeta_t,$$

where x_t is the (hidden) state, u_t is the input or control to the system, and y_t is the observation. The terms w_t, ζ_t are noise terms, and the matrices A, B, C, D are called the system matrices. A linear

¹This definition corresponds to the *valid* mode of convolution in typical implementations of convolution e.g. scipy.

216 dynamical predictor with parameters A, B, C, D predicts according to 217

 π

218

219

223

224

t-1

$$_{ABCD}(u_{1:t}, y_{1:t-1}) = \sum_{i=1} CA^{i-1}Bu_{t-i} + Du_t.$$

220 The best such predictor for a given sequence is also called the optimal open loop predictor, and it is 221 accurate if the signal is generated by a LDS without noise. 222

2.4 AUTO-REGRESSIVE SEQUENCE GENERATION FROM A PROMPT

Another mode of sequence prediction with large language models being its core use-case is that 225 of auto-regressive sequence generation starting from a prompt. Herein the sequence model has to 226 generate a specified number of tokens given a certain context. This is depicted in Figure 1. The 227 setting of auto-regressive generation from a prompt consists of two stages, the prefill stage and the 228 decode stage. During the prefill stage, the model ingests the context vector and generates a cache 229 that stores information required in the decode stage. 230

In the decode stage, the model takes the cache and the most recently generated token as input and 231 generates the next output token. Then the cache is updated with the most recent input token. We 232 denote the generation length at the decode stage with K. In contrast to pre-training, where the model 233 takes in a training sequence and predicts the next token, in the prefill generation setting the model 234 only has access to the cache and the most recent token when making a prediction. 235



240 241 242 243

244

245

251

256

260

236

237 238 239

Figure 1: Auto-regressive sequence generation from a prompt.

ABSTRACTING CONVOLUTIONAL SEQUENCE PREDICTION 2.5

246 We define a convolutional sequence prediction model to be given by a *filter*, which is a vector 247 denoted by $\phi \in \mathbb{R}^L$ where L is considered the *context length* of the model. It takes as an input a 248 sequence u, and outputs a prediction sequence. The above definition can be extended to multiple filter *channels* and nonlinearities, as we elaborate below with different examples. Formally, a single 249 output in the predicted sequence using a convolutional sequence model is given by 250

$$\hat{y}_t = \langle \phi, u_{t:t-L} \rangle. \tag{2}$$

252 This paradigm captures several prominent convolutional sequence models considered in the litera-253 ture. We highlight some of them below. The online convolution technique proposed by us can be 254 used with all the models below in straightforward manner leading to generation time improvement from $O(L^2)$ to $O(L \log^2 L)$. 255

State Space Models Discrete state space models such as those considered in Gu et al. (2021a) 257 have shown considerable success/adoption for long range sequence modelling. A typical state space 258 model can be defined via the following definition of a Linear Dynamical System (LDS) 259

$$x_t = Ax_{t-1} + Bu_t, y_t = Cx_t + Du_t$$
(3)

261 where u, y are the input and output sequences and A, B, C, D are the learned parameters. Various papers deal with specifications of this model including prescriptions for initialization (Gu et al., 262 2020), diagonal versions (Gupta et al., 2022), gating (Mehta et al., 2023) and other effective sim-263 plifications (Smith et al., 2023). All these models can be captured via a convolutional model by 264 noticing that the output sequence y in (3) can be written as 265

266
$$y = \phi * u + Du$$

267 where the filter ϕ takes the form $\phi_i = CA^{i-1}B$. Thus a convolutional sequential model.

odel with learnable filters ϕ generalizes these state space models. However, SSM are more efficient for generation and require only constant time for generating a token, where the constant depends on the size of the SSM representation.

LongConv/SGConv. The LongConv (Fu et al., 2023) and SGConv (Li et al., 2022) architectures, exploit the above connection and propose direct regularizations of the kernel to bias them towards kernels representing a state space model.

Spectral Transform Units. The STU architecture was proposed in Agarwal et al. (2023) based on the technique of spectral filtering for linear dynamical systems (Hazan et al., 2017; 2018). These are basically convolutional sequence models based on carefully constructed filters that are **not data dependent**. Rather, let $\phi_1, ..., \phi_k$ be the first k eigenvectors of the Hankel matrix H_L given by

$$H_L = \int_0^1 \mu_\alpha \mu_\alpha^\top d\alpha \in \mathbb{R}^{L \times L}, \ \mu_\alpha = (\alpha - 1)[1, \alpha, \alpha^2, ..., \alpha^{L-1}].$$

Then the STU outputs a prediction according to the following rule $\hat{y}_t = \sum_{i=1}^k M_i \langle \phi_i, u_{t:t-L} \rangle$, where ϕ_i are the eigenvectors as above and $M_{1:k}$ are learned projection matrices. The STU ar-chitecture is particularly appealing for learning in dynamical systems with long context, as it has theoretical guarantees for this setting, as spelled out in Agarwal et al. (2023).

Hyena. The Hyena architecture proposed in Poli et al. (2023), sequentially applies convolutions and element-wise products alternately. Formally, given an input $u_{1:T}$, N+1 linear projections $v, x_1, \ldots x_N$ of the input are constructed (similar to the q, k, v sequence in self-attention). The hyen operator as a sequence of convolution with learnable filters $h_1 \dots h_N$ is then given by

$$y = x^{N} \cdot \left(h^{N} * \left(x^{N-1} \cdot \left(h^{N-1} * (\ldots)\right)\right)\right)$$

EFFICIENT ONLINE CONVOLUTIONS USING FUTUREFILL

We begin by introducing a simple convenient primitive which we call FutureFill, which forms the crucial building block of our algorithms. Intuitively FutureFill corresponds to computing the *contribution* of the current and prevolusly generated tokens on the future tokens yet to be generated. For a convolutional model (and unlike attention) this contribution can be efficiently determined without even having generated the future tokens. Here onwards, for brevity of notation for any $v \in \mathbb{R}^{t}$, we assume $v_j = 0$ for any $j \leq 0$ or any j > t. Formally, given two sequences $v \in \mathbb{R}^{t_1}$, $w \in \mathbb{R}^{t_2}$ we define FutureFill $(v, w) \in \mathbb{R}^{t_2 - 1}$ as ³

$$\forall s \in [t_2 - 1] \; [\text{FutureFill}(v, w)]_s = \sum_{i=1}^{t_2 - s} v_{t_1 - i + 1} \cdot w_{s+i}$$



Figure 2: FutureFill operation between an input sequence and a convolutional filter.

Figure 2 depicts the FutureFill operation between an input sequence and a convolutional filter. Con-ceptually, [FutureFill(v, w)]_s is the contribution of the input v of length t_1 to the output [v * w] at position $t_1 + s$. The FFT algorithm for convolutions can easily be extended to compute the Future-Fill as well in time at most $O((t_1 + t_2) \log(t_1 + t_2))$. For example the *full* mode of a standard conv

²more precisely, there are additional linear and constant terms depending on the exact filters used, such as $\hat{y}_t = \hat{y}_{t-2} + \sum_{i=1}^3 M_i^u u_{t+1-i} + \sum_{i=1}^k M_i \langle \phi_i, u_{t:t-L} \rangle, \text{ see Agarwal et al. (2023) for more details.}$ ³recall that we denote $[x] = \{1 \dots x\}.$

implementation (e.g. scipy) can be used to compute FutureFill in the following way under Python
 slicing convention (exclusive of the last index)

FutureFill(v, w) = scipy.linalg.conv(v, w, mode=full)[t_1:t_1+t_2-1]

To leverage *FutureFill* into an efficient way to generate tokens from a convolutional model, consider the following simple proposition that follows from the definition of convolution.

Proposition 1. Given two vectors $a, b \in \mathbb{R}^t$, we have that

327

328

330

336

337 338

339 340

341

342

343 344

$$\forall t_1, s \in [t] \qquad [a * b]_s = \begin{cases} [a_{1:t_1} * b_{1:t_1}]_s & \text{if } s \leq t_1 \\ [a_{t_1+1:t} * b_{1:t_-t_1}]_{s-t_1} + [\text{FutureFill}(a_{1:t_1}, b)]_{s-t_1} & \text{otherwise} \end{cases}$$

We provide a proof of the proposition in the appendix. We use the above proposition to design efficient algorithms for online convolution.

3.1 EPOCHED-FUTUREFILL: EFFICIENT ONLINE CONVOLUTIONAL PREDICTION

When computing online convolutions, the FutureFill routine allows for the efficient pre-computation for the effect of past tokens on future tokens. We leverage this property towards online convolution via the Epoched-FutureFill procedure outlined in Algorithm 1.

Algorithm 1 Epoched-FutureFill: Efficient Online Convolutional Prediction

345 1: Input: Convolutional filter $\phi \in \mathbb{R}^L$. Input sequence $u \in \mathbb{R}^L$, streaming one coordinate every 346 round. K, the epoch length. 347 2: Set $\tau = 1$. Set FutureFill cache $C \in \mathbb{R}^K$ to 0. 348 3: for t = 1, 2, ..., L do 349 4: Receive u_t . Compute and output $\hat{y}_t = \sum_{j=1}^{\tau} u_{t+1-j} \cdot \phi_j + C_{\tau}$ 350 5: 351 6: if $\tau = K$ then 352 Compute FutureFill cache $C \in \mathbb{R}^K$ defined as $C_i = [\text{FutureFill}(u_{1:t}, \phi_{1:t+K})]_i$. 7: 353 8: $\tau \leftarrow 1$ 354 9: else 10: $\tau \leftarrow \tau + 1$ 355 end if 11: 356 12: end for 357 358 359 360 361 362 FutureFill cost – $O(L \log L)$ Compute k^{th} token - O(k)Overall: $\frac{\sum_{l=1}^{K} i + L \log L}{\nu} = \frac{K^2 + L \log L}{\nu} \sim O(\sqrt{L \log L})$ 364 365 366 Cache Size K 367 368 Figure 3: Illustration for Algorithm 1 369 370 In the following lemma we state and prove the properties that Epoched-FutureFill enjoys. The 371 theorem provides a trade-off between the additional memory overhead and total runtime incurred 372 by the algorithm. In particular, the runtime in this tradeoff is optimized when the total memory is $O(\sqrt{L \log L})$ leading to a total runtime of $O(L^{3/2}\sqrt{\log L})$. 373 374 **Theorem 2.** Algorithm 1 computes the online convolution of sequences with length L and runs in 375 total time $O\left(\frac{L^2 \log L}{K} + KL\right)$ with a total additional memory requirement of O(K). In particular 376 setting $K = \sqrt{L \log L}$, we get that Algorithm 1 computes online convolution in $O(L^{3/2}\sqrt{\log L})$ 377 total time and $O(\sqrt{L \log L})$ memory.

378 *Proof.* Since the proof of correctness is mainly careful accounting the contributions for various 379 indices, we provide it in the appendix. We prove the running time bounds below. The running time 380 consists of two components as follows: 381

- 1. Every iteration, line 5 is executed. One term, C_{τ} , has already been computed and saved in line 7. We can retrieve it in time O(1). The other term is a sum of τ products, which can be computed in time τ .
- 2. Every K iterations, we execute line 7 and update the terms in the cache. The FutureFill operation can be computed via the FFT taking at most $O(L \log L)$ time.

The overall running time is computed by summing over the L iterations. In each block of K iterations, we apply FFT exactly once, and hence the total computational complexity is

$$\frac{L}{K} \left(L \log L + \sum_{\tau=1}^{K} \tau \right) = O\left(\frac{L^2 \log L}{K} + KL \right) = O\left(L^{3/2} \sqrt{\log L} \right),$$

where the last equality holds when the cache size K is chosen to minimize the sum, i.e. K = $\sqrt{L\log L}$.

CONTINUOUS-FUTUREFILL: QUASILINEAR ONLINE CONVOLUTIONAL PREDICTION 3.2

In this section we provide a procedure that significantly improves upon the runtime of Epoched-400 FutureFill. Our starting point is Proposition 1, which implies that, to compute the convolution 401 between two sequences we can break the sequences at any point, compute the convolution between 402 the corresponding parts and *stitch* them together via a FutureFill computation. This motivates the following Divide and Conquer algorithm to compute the convolution of two sequences $a, b \in \mathbb{R}^{L}$

- Recursively compute $a_{1:L/2} * b_{1:L/2}, a_{L/2+1:t} * b_{1:L/2}$.
- 406 407

403

404 405

382

384

385 386

387 388 389

395

396 397

398 399

• Output the concatenation of $a_{1:L/2} * b_{1:L/2}$ and $(a_{L/2+1:t} * b_{1:L/2}) + \text{FutureFill}(a_{1:L/2}, b)$.

408 Since FutureFill for L length sequences can be computed in time $O(L \log L)$ via the FFT, it can be 409 seen via the standard complexity calculation for a divide and conquer algorithm that the computa-410 tional complexity of the above algorithm in total is $O(L \log^2 L)$. As an offline algorithm, this is 411 naturally worse than the computational complexity of FFT itself, however as we show in the fol-412 lowing, the advantage of the above algorithm is that it can be executed in an *online* fashion, i.e. the 413 tokens can be generated as the input streams in, with the same computational complexity.

414 We provide a formal description of the algorithm in Algorithm 2. We note that the formal description 415 of the above algorithm essentially serializes the sequence of operations involved in the above divide 416 and conquer procedure by their chronological order. For high-level intuition we encourage the reader 417 to maintain the divide and conquer structure when understanding the algorithm. The algorithm 418 proceeds as follows: at each time step, $\hat{y}_t = \langle u_{1:t}, \phi_{t:1} \rangle$ is returned as a sum of C_t , the cache that stores the contribution from past tokens, and $u_t \cdot \phi_1$, the contribution from token u_t . In Line 7, 419 the algorithm then computes the contribution of tokens $u_{t-2^{k(t)}+1:t}$ to positions $t+1, \ldots, t+2^{k(t)}$ 420 421 of $[u * \phi]$. Finally, we add the output of FutureFill to the existing cache C to accumulate the contributions. In Figure 4, we provide an execution flow for the algorithm for convolving two 422 sequences of length 8 highlighting each FutureFill operation that is computed. 423

424 In the following theorem we prove a running time bound for Algorithm 2. We provide the proof of 425 correctness in the appendix, as it boils down to accounting of contribution from various parts.

426 **Theorem 3.** Algorithm 2 computes the online convolution of sequences with length L and runs in 427 total time $O(L \log^2(L))$ with a total additional memory requirement of O(L). 428

429

Proof. As can be seen from the algorithm for every generated token the most expensive operation is 430 the FutureFill computed in Line 6 so we bound the total runtime of that operation. Note that at any 431 time t, the cost of FutureFill operation is $O((1 \vee k(t)) \cdot 2^{k(t)})$, where $a \vee b$ denotes the max of a and



Figure 4: Quasilinear Online Convolution using FutureFill: Figure shows the execution flow for Algorithm 2 for convolving 8-length sequences. The input sequence u streams in an online fashion and the filter ϕ is fully available to the algorithm. The colors are representative of the size of the FutureFill operations performed and the time t (also appropriately color-coded) highlights when the FutureFill operations were performed.

b. Summing this over every time step t we get,

$$\sum_{t=1}^{L} (1 \lor k(t)) 2^{k(t)} = \sum_{k=0}^{\lfloor \log L \rfloor} |\{t : k(t) = k\}| (1 \lor k) 2^k \le L + \sum_{k=1}^{\lfloor \log L \rfloor} 2^{\lfloor \log L \rfloor - k + 1} \cdot k 2^k \le 3L \sum_{k=1}^{\lfloor \log L \rfloor} k \le 3L \log^2 L.$$

Thus the total runtime of the algorithm is bounded by $O(L \log^2 L)$.

Algorithm 2 Continuous-FutureFill: Quasilinear Generation From Convolutional Models

1: Input: Convolutional filter $\phi \in \mathbb{R}^L$. Input sequence $u \in \mathbb{R}^L$, streaming one coordinate every round. 2: Set $b = \lfloor \log L \rfloor$. Set FutureFill cache $C \in \mathbb{R}^L$ to 0.

3: for t = 1 ... L do

Receive u_t . Output $\hat{y}_t = C_t + u_t \cdot \phi_1$. 4:

5: Let k(t) be the highest power of 2 that divides t, i.e. $k = \max\{i \in [b] : t \mod 2^i = 0\}$. Compute FF = FutureFill($u_{t-2^{k(t)}+1:t}, \phi_{1:2^{k(t)}+1}$) 6: Set $C_i = C_i + FF_{i-t}$ $\forall i \in [t+1, t+2^{k(t)}]$ 7: 8: end for

9:

FAST AUTO-REGRESSIVE SEQUENCE GENERATION FROM A PROMPT

In this section we consider the problem setting of auto-regressively generating K tokens starting from a given prompt of length L. For convolutional models specifically we define an abstract version of the problem as follows, given a prompt vector $p \in \mathbb{R}^L$ and a convolutional filter $\phi \in \mathbb{R}^{L+K}$ ⁴, the aim is to iteratively generate the following sequence of tokens

$$\hat{y}_t = \langle \hat{y}_{1:t-1}, \phi_{t-1:1} \rangle + \langle p_{1:L}, \phi_{t+L-1:t} \rangle = \sum_{j=1}^{t-1} \hat{y}_{t-j} \cdot \phi_j + \sum_{j=t}^{t+L-1} p_{t+L-j} \phi_j.$$

⁴the assumption of the filter being larger than L + K is without loss of generality as it can be padded with 0s

As can be seen from the above definition the expected output is an online convolution where the input sequence u has a prefix of the prompt p and the input sequence is appended by the most recently generated output by the model (i.e. auto-regressive generation). We note that we only consider the *convolution* part of a convolutional model (eg. STU) above for brevity and other parts like further projection of the tokens etc can be appropriately added. As mentioned the above model naturally fits into online convolution and the following algorithm delineates the method to use ContinuousFutureFill (Algorithm 2) for the above problem.

Algorithm 3 Fast auto-regressive sequence generation from a prompt using FutureFill

1: **Input:** K > 0, L > 0, prompt $p_{1:L}$, convolutional filter $\phi \in \mathbb{R}^{L+K}$.

- 2: Set up a FutureFill cache $C \in \mathbb{R}^{K}$ as $C \leftarrow$ FutureFill (p, ϕ) .
- 3: Set up the online convolution algorithm (Algorithm 6) with filter ϕ and sequence length K, i.e. $\mathcal{A} \leftarrow \text{ContinuousFutureFill}(\phi)$.

4: Running candidate token $y \leftarrow 0$.

5: for t = 1, ..., K do

6: Output $\hat{y}_t \leftarrow C_t + y$.

7: Generate next token candidate $y \leftarrow \mathcal{A}(\hat{y}_t)$.

8: end for

493

494

495

496

497

498

499

500

501

502

504

505

506

507

508 509

510 511

512

513

514

515

516

517

518

519

532

533

534 535 The correctness of the algorithm is immediate via the properties of FutureFill and ContinousFuture-Fill. The following corollary bounding running time also follows easily from Theorem 3.

Corollary 4. Algorithm 3 when supplied with a prompt of sequence length L, generates K tokens in total time $O(L \log L + K \log^2 (L + K))$ using a total cache of size O(K).

5 EXPERIMENTS

In this section, we use a convolutional model that generates tokens in an online fashion to verify our results. We experimentally evaluate Epoched-FutureFill (Algorithm 1) which has a runtime of $O(L^{3/2}\sqrt{\log L})$ and Continuous-FutureFill (Algorithm 2) which has a runtime of $O(L \log^2 L)$ against the naive implementation which has a runtime of $O(L^2)$ when generating L tokens from scratch. For increasing values of L, we measure the time S(L) it takes for a single layer to generate L tokens. In Figures 5 and 6 we plot the amortized step time S(L)/L and total generation time S(L), respectively, as functions of L. We see the behavior that is expected: the naive decoder runs in amortized O(L) per step, while our methods achieve sublinear and logarithmic decoding complexities respectively.



Figure 5: Average number of seconds per step when generating L tokens, as a function of L.



Figure 6: Total number of seconds to generate L tokens, as a function of L.

Due to differences in hardware acceleration, inference pipeline implementation, and other engineering details, it would be difficult to present timing results with a properly-optimized setup. On large
decoding platforms involving prefill caching, these variations only become more complicated. We
opted to time things for one layer on CPU in a simple online decoding loop with a large number of
tokens to make the asymptotic gains clear.

540 REFERENCES

556

586

592

- Naman Agarwal, Daniel Suo, Xinyi Chen, and Elad Hazan. Spectral state space models. arXiv preprint arXiv:2312.06837, 2023.
- James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao Zhang. JAX: composable transformations of Python+NumPy programs, 2018. URL http: //github.com/jax-ml/jax.
- Tri Dao, Daniel Y Fu, Khaled K Saab, Armin W Thomas, Atri Rudra, and Christopher Ré.
 Hungry hungry hippos: Towards language modeling with state space models. *arXiv preprint arXiv:2212.14052*, 2022.
- Dan Fu, Simran Arora, Jessica Grogan, Isys Johnson, Evan Sabri Eyuboglu, Armin Thomas, Benjamin Spector, Michael Poli, Atri Rudra, and Christopher Ré. Monarch mixer: A simple subquadratic gemm-based architecture. *Advances in Neural Information Processing Systems*, 36, 2024.
- Daniel Y Fu, Elliot L Epstein, Eric Nguyen, Armin W Thomas, Michael Zhang, Tri Dao, Atri Rudra, and Christopher Ré. Simple hardware-efficient long convolutions for sequence modeling. *arXiv* preprint arXiv:2302.06646, 2023.
- Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv preprint arXiv:2312.00752*, 2023.
- Albert Gu, Tri Dao, Stefano Ermon, Atri Rudra, and Christopher Ré. Hippo: Recurrent memory with optimal polynomial projections. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 1474–1487. Curran Associates, Inc., 2020.
- Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with structured state spaces. *arXiv preprint arXiv:2111.00396*, 2021a.
- Albert Gu, Isys Johnson, Karan Goel, Khaled Saab, Tri Dao, Atri Rudra, and Christopher Ré. Combining recurrent, convolutional, and continuous-time models with linear state space layers. Advances in neural information processing systems, 34:572–585, 2021b.
- Ankit Gupta, Albert Gu, and Jonathan Berant. Diagonal state spaces are as effective as structured state spaces. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), Advances in Neural Information Processing Systems, 2022. URL https://openreview.net/forum?id=RjS0j6tsSrf.
- 578 Elad Hazan and Karan Singh. Introduction to online nonstochastic control. *arXiv preprint* 579 *arXiv:2211.09619*, 2022.
- Elad Hazan, Karan Singh, and Cyril Zhang. Learning linear dynamical systems via spectral filtering. In Advances in Neural Information Processing Systems, pp. 6702–6712, 2017.
- Elad Hazan, Holden Lee, Karan Singh, Cyril Zhang, and Yi Zhang. Spectral filtering for general linear dynamical systems. In *Advances in Neural Information Processing Systems*, pp. 4634–4643, 2018.
- Elad Hazan et al. Introduction to online convex optimization. Foundations and Trends[®] in Optimization, 2(3-4):157–325, 2016.
- ⁵⁸⁹ Norman P Jouppi, Doe Hyun Yoon, George Kurian, Sheng Li, Nishant Patil, James Laudon, Cliff
 ⁵⁹⁰ Young, and David Patterson. A domain-specific supercomputer for training deep neural networks.
 ⁵⁹¹ *Communications of the ACM*, 63(7):67–78, 2020.
- 593 Yuhong Li, Tianle Cai, Yi Zhang, Deming Chen, and Debadeepta Dey. What makes convolutional models great on long sequence modeling? *arXiv preprint arXiv:2210.09298*, 2022.

- Y Isabel Liu, Windsor Nguyen, Yagiz Devre, Evan Dogariu, Anirudha Majumdar, and Elad Hazan.
 Flash stu: Fast spectral transform units. *arXiv preprint arXiv:2409.10489*, 2024.
- Stefano Massaroli, Michael Poli, Dan Fu, Hermann Kumbong, Rom Parnichkun, David Romero,
 Aman Timalsina, Quinn McIntyre, Beidi Chen, Atri Rudra, et al. Laughing hyena distillery:
 Extracting compact recurrences from convolutions. *Advances in Neural Information Processing Systems*, 36, 2024.
- Harsh Mehta, Ankit Gupta, Ashok Cutkosky, and Behnam Neyshabur. Long range language model ing via gated state spaces. In *The Eleventh International Conference on Learning Representations*, 2023.
- Antonio Orvieto, Samuel L Smith, Albert Gu, Anushan Fernando, Caglar Gulcehre, Razvan Pascanu, and Soham De. Resurrecting recurrent neural networks for long sequences. *arXiv preprint arXiv:2303.06349*, 2023.
- Michael Poli, Stefano Massaroli, Eric Nguyen, Daniel Y Fu, Tri Dao, Stephen Baccus, Yoshua
 Bengio, Stefano Ermon, and Christopher Ré. Hyena hierarchy: Towards larger convolutional
 language models. In *International Conference on Machine Learning*, pp. 28043–28078. PMLR, 2023.
- ⁶¹² Noam Shazeer. Glu variants improve transformer. *arXiv preprint arXiv:2002.05202*, 2020.
- Jiaxin Shi, Ke Alexander Wang, and Emily Fox. Sequence modeling with multiresolution convolutional memory. In *International Conference on Machine Learning*, pp. 31312–31327. PMLR, 2023a.
- Jiaxin Shi, Ke Alexander Wang, and Emily B. Fox. Sequence modeling with multiresolution convolutional memory, 2023b. URL https://arxiv.org/abs/2305.01638.
- Jimmy T.H. Smith, Andrew Warrington, and Scott Linderman. Simplified state space layers for sequence modeling. In *The Eleventh International Conference on Learning Representations*, 2023.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. Advances in neural information processing systems, 30, 2017.

⁶⁴⁸ A UPDATED EXPERIMENTS SECTION AND DETAILS

649

663

665

666

667

668

669

670

671 672

673 674

675

676

677

650 In this section, we employ a convolutional sequence prediction model that generates tokens in an 651 online fashion to verify our results. We experimentally evaluate Epoched-FutureFill (Algorithm 652 1) which has a runtime of $O(K^{3/2}\sqrt{\log K})$ and Continuous-FutureFill (Algorithm 2) which has a 653 runtime of $O(K \log^2 K)$ against the naive implementation of convolution which has a runtime of 654 $O(K^2)$ when generating K tokens from scratch. We also provide a comparison with a self-attention 655 based Transformer model (with a standard implementation of KV cache and with the same hidden 656 dimension, number of layers and commensurately chosen other parameters, see next subsection for complete details on these models) 657

For increasing values of K, we measure the time it takes for the model to generate K tokens from scratch (i.e. no prompt provided). In Figure 7 we plot the amortized step time the total generation time, as functions of K. We see the behavior that is expected: the naive decoder runs in total time $O(K^2)$ per step, similar to the decoder for transformer while our method EpochedFutureFill is able to achieve a significant sub-quadratic improvement.



Figure 7: Total time for generating K tokens, as a function of K.

678 679 680

681

697

698

A.1 EXPERIMENT DETAILS

682 For our experiments we consider a two layer model with either multi-headed self-attention layers 683 (referred to as Transformer) or the STU layers (referred to as convolutional network). The hidden 684 dimension (d) of the network is fixed to be 32 with the number of heads fixed to be 4 and the 685 key/value size kept at 8. The networks contain standard implementations of residual connections, layer-norms and a feed-forward (FFN) layer in between every layer. The FFN layer used in the 686 experiments is the FFN_{GeGLU} layer proposed in Shazeer (2020). For the attention layers we employ 687 a standard implementation of KV-cache for efficiency (i.e. caching the KV values for previously 688 generated token for every layer). 689

For the STU layer we use the tensored-approximation experimented upon in Liu et al. (2024). Note that during inference for the tensored-approximation of STU, the layer maintains as parameters two matrices $M_{input} \in \mathbb{R}^{d \times d}$ and $M_{filters} \in \mathbb{R}^{(K) \times d}$, where K is the number of generated tokens and d is the hidden dimensionality. We provide a detailed equation describing the exact operation performed by the STU layer in generating the k^{th} token below. Let $x_1 \dots x_K$ be the embeddings of tokens generated in online manner, i.e. when generating the k^{th} token only the embeddings $x_1 \dots x_{k-1}$ are available to the model. The generated token sequence follows the following implicit equation

$$[x_k \dots x_1] = \left[M_{\text{filters}} * \left(M_{\text{inputs}} [x_{k-1} \dots x_1] \right) \right]_{1:k}.$$

The above convolution operation is over d-dimensional sequences which is implemented as d 1dimensional convolutions performed along each dimension. Since we have equated the hidden dimensionality of the network across all our settings, we can see that, naively computed, the number of flops per token of both the self-attention model as well as the convolutional model are of the same order which is also observed in the experiments. The next section provides details on how to
 leverage FutureFill based online convolution algorithms for the STU layer which we implement in
 our experiments.

Finally all our experiments are implemented in Jax Bradbury et al. (2018) were performed on a single Google TPUv2 machine (Jouppi et al. (2020)).

708 709

710

711

712

713 714

715

740

741 742

743 744 745

746 747

754 755

B FAST ONLINE CONVOLUTIONAL PREDICTION

The techniques for online convolution can naturally be applied to online prediction using a convolutional model. We demonstrate this use case via its application to the STU architecture proposed in Agarwal et al. (2023) based on the spectral filtering algorithm (Hazan et al., 2017).

716 717 We illustrate in more detail how the method works for the STU model in Algorithm 4. It improves 718 the total running time from $O(L^2)$ of the original spectral filtering algorithm from Hazan et al. 719 (2017) to $O(L \log^2 L)$ while maintaining the same regret bound.

| 1: | Input: $K > 0, L > 0.$ |
|-----|---------------------------------------------------------------------------------------------------------------------------------------------------------|
| 2: | Set variables $\{M_1^1 \dots M_K^1 \in \mathbb{R}^{d_{out} \times d_{in}} \leftarrow 0\}$ and set $\{\phi_1 \dots \phi_K\}$ as the largest eigenvectors |
| | of H_L , the Hankel matrix corresponding to length-L sequences. |
| 3: | Initialize K OnlineConvolution modules, one for each filter $\{\mathcal{A}_k(\phi_k)\}_{k=1}^K$. |
| 4: | for $t = 1, 2,, L$ do |
| 5: | Receive input token u_t . |
| 6: | for $k = 1, 2, \dots K$ do |
| 7: | $F_k \leftarrow \mathcal{A}_k(\phi_k)(u_t).$ |
| 8: | end for |
| 9: | Compute and predict $\hat{y}_t = \sum_{k=1}^K M_k^t F_k$. |
| 10: | Observe y_t , suffer loss $\ell_t(M_{1,k}^t) = y_t - \hat{y}_t ^2$, and update $M_{1,k}^{t+1} \leftarrow \nabla \ell_t(M_{1,k}^t)$. |
| 11: | end for |

The main claim regarding the performance of Algorithm 4 follows directly from Theorems 2 and 3 and is as follows.
The main claim regarding the performance of Algorithm 4 follows directly from Theorems 2 and 3 and is as follows.

Corollary 5. Algorithm 4 with sequence length L guarantees the same regret bound as spectral filtering (Hazan et al., 2017) with context length L. Furthermore its computational complexity based on the online convolution module used are as follows:

• If using EpochedFutureFill(Algorithm 1): Runtime - $O(L^{3/2}\sqrt{\log L})$, Memory - $O(\sqrt{L \log L})$.

• If using ContinuousFutureFill(Algorithm 2): Runtime - $O(L \log^2 L)$, Memory - O(L).

C MISSING PROOFS

Proof of Proposition 1. Note that by definition, $[a * b]_s = \sum_{i=1}^s a_i b_{s+1-i}$. We now consider the two cases: for $s \leq t_1$, we have that

$$[a_{1:t_1} * b_{1:t_1}]_s = \sum_{i=1}^s a_i b_{s+1-i} = [a * b]_s$$

For the case when $t \ge s > t_1$, we have that

$$[a_{t_1+1:t} * b_{1:t-t_1}]_{s-t_1} = \sum_{i=1}^{s-t_1} a_{t_1+i} b_{s-t_1+1-i} = \sum_{i=t_1+1}^{s} a_i b_{s+1-i},$$

where the last equality follows by redefining $i = t_1 + i$. Further we have that

FutureFill
$$(a_{1:t_1}, b)$$
] _{$s-t_1$} = $\sum_{i=1}^{t-s+t_1} a_{t_1-i+1} \cdot b_{s-t_1+i} = \sum_{i=1}^{t_1} a_{t_1-i+1} \cdot b_{s-t_1+i} = \sum_{i=1}^{t_1} a_i \cdot b_{s+1-i}$,
For $a_{t_1-i+1} \cdot b_{s-t_1+i} = \sum_{i=1}^{t_1} a_{t_1-i+1} \cdot b_{s-t_1+i} = \sum_{i=1}^{t_1} a_{t_1-i+1} \cdot b_{s-t_1+i} = \sum_{i=1}^{t_1} a_{t_1-i+1} \cdot b_{s-t_1+i}$

where the second last equality follows by noting that a_i is assumed to be 0 for all $j \leq 0$ and the last equality follows by redefining $i = t_1 - i + 1$. Overall putting the two together we get that

763
764
$$[a_{t_1+1:t}*b_{1:t-t_1}]_{s-t_1} + [FutureFill(a_{1:t_1}, b)]_{s-t_1} = \sum_{i=1}^{t_1} a_i \cdot b_{s+1-i} + \sum_{i=1}^{t_1} a_i \cdot b_{s+1-i} = \sum_{i=1}^{s} a_i \cdot b_{s+1-i} = [a*b]_s$$

765
766 This finishes the proof.

This finishes the proof.

Proof of correctness for Algorithm 1. Consider any time t and the output \hat{y}_t . Let $t' \leq t$ be the last time when Line 7 was executed, i.e. FutureFill was computed. By definition $t' = t - \tau$. Note the following computations.

$$\hat{y}_t = \sum_{j=1}^{\tau} u_{t+1-j} \cdot \phi_j + C_{\tau} = \sum_{j=1}^{\tau} u_{t+1-j} \cdot \phi_j + [\text{FutureFill}(u_{1:t'}, \phi_{1:t'+K})]_{\tau}$$

 $=\sum_{i=1}^{\tau} u_{t+1-j} \cdot \phi_j + \sum_{i=1}^{t'+K-\tau} u_{t'-j+1} \cdot \phi_{\tau+j}$

778
779
$$= \sum_{j=1}^{\tau} u_{t+1-j} \cdot \phi_j + \sum_{j=1}^{t'} u_{t'-j+1} \cdot \phi_{\tau+j}$$
780

781
782
$$-\sum_{i=1}^{\tau} u_{i+1} \cdots v_{i+1} \sum_{i=1}^{t-\tau} u_{i+1} \cdots u_{i+1}$$

$$= \sum_{j=1}^{r_{02}} u_{t+1-j} \cdot \phi_j + \sum_{j=1}^{r_{02}} u_{t-\tau-j+1} \cdot \phi_{\tau+j}$$
783
784
 τ
 t

784
785
$$= \sum_{j=1}^{\tau} u_{t+1-j} \cdot \phi_j + \sum_{j=\tau+1}^{t} u_{t-j+1} \cdot \phi_j = [u * \phi]_t$$
786

| - | | | - | |
|---|---|---|---|---|
| 7 | s | 2 | - | 7 |
| 1 | c | | 1 | |
| | | | | |
| | | | | |

| 7 | 8 | 8 |
|---|---|---|
| 7 | 8 | 9 |

Proof of correcteness for Algorithm 2. We will focus on showing that $C_t = \sum_{i=2}^t u_{t+1-i}\phi_i$. Since the output is $C_t + u_t \cdot \phi_1$, this will suffice for the proof. For brevity of the proof and without loss of generality we will assume L is a power of 2. For cleaner presentation for the s^{th} coordinate of vector v we will use the notation v_s and v|s| interchangeably in this section.

We first introduce some definitions for convenience in this section. Given an index $i \leq L$ we define its decomposition $\{i_1, i_2 \dots i_m\}$ as the unique sequence of numbers $\leq \log L$ such that following holds

$$i_1 > i_2 > i_3 \dots$$
 and $i = \sum_j 2^{i_j}$.

These indices correspond to the ones in a log L-bit representation of i. Note that k(i) as defined in the algorithm is equal to i_m . Further we define the cumulants of i as the following sequence of numbers $\{i'_1, i'_2 \dots\}$ satisfying

$$i_{\tau}' = \sum_{j=1}^{\tau} 2^{i_j}.$$

Thus we have that $i'_1 < i'_2 < \ldots i'_m = i$. We now prove the following lemma which specifies when the FutureFill cache gets updated in an execution of the algorithm.

Lemma 6. Given an index $i \leq L$, consider its decomposition $\{i_1, i_2 \dots i_m\}$ and cumulants $\{i'_1, i'_2 \dots i'_m\}$ as defined above. It holds that the value of C_{i+1} is updated (as in Line 8 in the algorithm) only when t is one of $\{i'_1, i'_2 \dots i'_m\}$.

A direct consequence of the above lemma is that given any index i we have that the value of C_{i+1} is not updated after time step i. Further using the decomposition $\{i_1, i_2 \dots i_m\}$ and cumulants $\{i'_1, i'_2 \dots i'_m\}$ of i and the update equations for C (Line 8), we have that final value of C_{i+1} is given by the following,

$$C_{i+1} = \sum_{j=1}^{m} \text{FutureFill}(u[i'_j - 2^{i_j} + 1 : i'_j], \phi[1 : 2^{i_j+1}])[i+1 - i'_j]$$

$$= \sum_{j=1}^{m} \sum_{k=1}^{2^{i_j}} u[i'_j - k + 1] \cdot \phi[i + 1 - i'_j + k]$$

820
821
822
823
824

$$j=1 \ k=1$$

 i'_j
 $r=i'_j-2^{i_j}+1$
 $u[r] \cdot \phi[i+1-r+1]$

$$= \sum_{r=1}^{i} u[r] \cdot \phi[i+1-r+1]$$

Thus the output of the algorithm for any *i*, satisfies

$$\hat{y}_{i+1} = C_{i+1} + u_{i+1} \cdot \phi_1 = \sum_{r=1}^{i} u[r] \cdot \phi[i+1-r+1] + u_{i+1} \cdot \phi_1 = \sum_{r=1}^{i+1} u[r] \cdot \phi[i+1-r+1] = [u*\phi]_{i+1}.$$

This proves the requisite. We finally provide a proof of Lemma 6 to finish the proof.

Proof of Lemma 6. By the definition of the algorithm, to be able to update C_{i+1} at some time t < i + 1 it must be the case that

$$i+1 \in [t+1, t+2^{k(t)}].$$

Consider some t and its decomposition $\{t_1, t_2 \dots t_n\}$ and cumulants $\{t'_1, t'_2 \dots t'_n\}$. By the definition of the update in Line 8, we have that at time t we only update indices i + 1 for which i has the sequence $\{t'_1, t'_2 \dots t'_{n-1}\}$ in its decomposition as a prefix. It can then be seen that for a given number i, the only such numbers are its cumulants, i.e. $\{i'_1 \dots i'_m\}$ which finishes the proof. \Box