

490 A More detailed comparisons with existing baselines

491 This section provides the reader with a more in-depth comparison with similar architectures. We
492 cover BRecT [20] in Section A.1 and GSS-Hybrid [24] in Section A.2

493 A.1 Comparison with Block Recurrent Transformer (BRecT)

494 The Block Transformer layer (i.e Slide:12L) also processes keys and values from the previous window
495 stored in a differentiable cache. This is implemented similarly to the sliding window attention pattern
496 suggested in [20] and was originally introduced by Transformer-XL [8]. Using a causal mask, at
497 every token inference step, the attention mechanism is applied to blocks of tokens of size W and is
498 partially extended to the cached keys and values from the previous block with the sliding window.
499 BRecT, as explained in [20], uses a non-differentiable cache that is carried from one sequence of size
500 L to the next². The last recurrent states of a sequence are stored in a non-differentiable cache and fed
501 to the next training step on the following sequence in the document as a warm-start. We do not pass
502 such a representation, since to compute the output of the convolution, we need access to the whole
503 sequence. We believe that this is an one advantage that BRecT has over our method, especially for
504 very long examples that split into ordered sequences of length L , since the cache carried from one
505 sequence to the next can provide very useful long range information and (weak) access to the whole
506 past. Since we need the whole sequence to compute SSM states, history beyond L may be lost in the
507 process. We believe that BST can further be improved by adding non-differentiable sequence cache
508 for very long documents.

509 While in other architectures, the history between blocks of tokens is not modeled, both BST and
510 BRecT use a mechanism to model previous block context. The authors of BRecT experiment with
511 various sequential gating mechanisms to condense the information from past blocks. With BST, we
512 use SSM to provide context from previous blocks to the current block as explained in Section 3.2

513 A.2 Comparison with the Transformer GSS-Hybrid

514 GSS-Hybrid [24] is a SSM-Transformer hybrid architecture that we first describe in Section 4.1. The
515 architecture is significantly different from BRT. GSS-Hybrid is primarily composed of Gated State
516 Space (GSS) layers and has a few interleaved Transformer layers at every 4th layer starting with the
517 2nd layer. BRT on the other hand is mainly composed of Block Transformer layers and has Block
518 State Transformer layers at layer positions {1, 7, 9} for the $\sim 200M$ model and {1, 5, 7, 9} for the
519 $\sim 400M$ model. Our hybrid does not stack SSM and Transformer layers like the GSS-Hybrid but rather
520 replaces the recurrence in BRecT with an SSM such as S4. In BRT, the SSM generates states for
521 each Block Transformer representations and we then use cross-attention to mix the states and the
522 self-attention outputs. We also use a simpler SSM. The authors in [24] initially built GSS, a gated
523 version of DSS [15], to (1) reduce SSM parameter dimensions, (2) stabilize training of the SSM and
524 (3) allow better length generalization. However, when experimenting with SSMs such as S4 or DSS,
525 we found that the gating was not necessary to achieve all three objectives stated above. We decided
526 that using GSS’s Gated Attention Unit [19] was therefore not needed when integrating SSM states
527 into the attention mechanism. We also reiterate that the authors in [24] used hyperparameter search
528 to get the best performance while we did not.

529 B Evaluating Length Generalization capabilities

530 We present our length generalization analysis and report perplexity in Figure 4. Our models and
531 baselines all have $\sim 400M$ parameters, are trained on a sequence length of 4k and tested on sequences
532 with *lower* and *higher* sequence lengths of {512, 16k, 65k}.

533 We notice that all models have similar perplexity for sequence lengths of 512. Both BST:SH:S4-L
534 and GSS-Hybrid-L generalize well on 16k and 65k sequence lengths for PG19 and arXiv. For GitHub,
535 GSS-Hybrid-L and BST:MF:unstruct-L perplexities increase drastically, potentially due to noise in the
536 GitHub dataset. For GitHub again, BRecT:fixed:skip-L performs very well at higher sequence lengths.
537 We hypothesize that the block recurrent model’s access to the entire past, via non-differentiable cache

²In our work and in [20], a document is split into multiple sequences of size L and each sequence is split into multiple blocks of size W

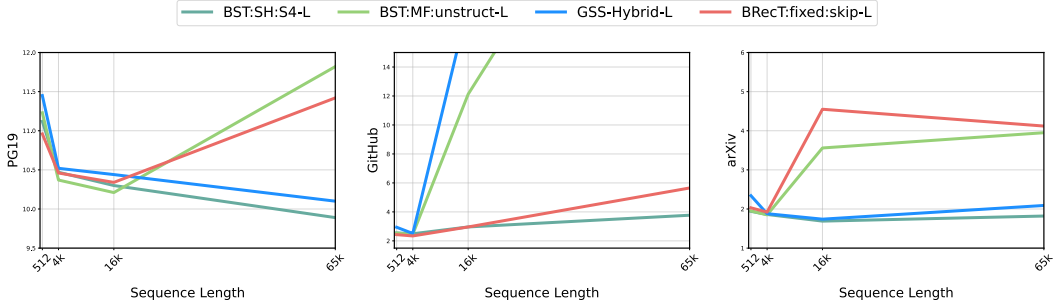


Figure 4: Length Generalization for sequence lengths $\{512, 16k, 65k\}$ on PG19 (left), GitHub (middle) and arXiv (right). BST:SH:S4-L generalizes better than any other baselines, including GSS-Hybrid-L that uses GSS, a structured SSM. GSS-Hybrid-L numbers are from [24].

538 of representations across sequences, helps retain a “memory” of dependencies between each code
 539 file in the GitHub dataset. Interestingly, we also note that BST:MF:unstruct-L and BRecT:fixed:skip-L
 540 outperform other methods on PG19 up to a sequence length of 16K. Perplexity performance on PG19
 541 is perhaps less reliant on long term relationships between tokens, which can explain the performance
 542 of models that have no explicit built-in mechanisms for length generalization.

543 The analysis also allows us to draw a clear distinction between *structured* and *unstructured* SSMs
 544 integrated in hybrid architectures. As previously mentioned in Section 3.1, SSMs such as GSS, DSS
 545 and S4 use a structured kernel K , built from learned matrices A , B and C for any sequence length L
 546 in Equation 3. Since K is extendable to any arbitrary sequence length L , both BST:SH:S4-L and GSS-
 547 Hybrid-L have a build-in mechanism for length generalization that the unstructured BST:MF:unstruct-L
 548 model does not. BST:MF:unstruct-L performs best on the training sequence of 4K and is on-par for
 549 512 with perplexity increasing for unseen 16K and 65K sequence lengths. **BST:SH:S4-L has by far**
 550 **the best perplexity for 65K sequence lengths on PG19, GitHub and arXiv.**

551 C Ablation Studies

552 In the following section, we perform ablations to investigate (1) the placement of a *single* SSM layer
 553 in Table 2 in the overall architecture, (2) the effects of the number of SSM layers added in Table 3
 554 and (3) the size D of the SSM state in Table 4. For the ablations, we use the $\sim 200M$ parameter
 555 BST:SH:S4, since it is the fastest model, and assess various configurations on PG19.

Table 2: A single BST at various layer index.

Layer index	Perplexity
3	12.41
7	11.92
9	11.88
12	12.03

Table 3: Multiple BST layers at various locations.

Num layers	Perplexity
2	11.69
3	11.57
4	11.21
5	11.20

Table 4: Increasing BST’s S4 model state size D .

State Size	Perplexity	Step Time
8	11.95	$\times 0.7$
16	11.57	$\times 1.0$
32	11.55	$\times 1.8$
64	11.54	$\times 3.2$

556 In Table 2, we experiment adding a single BST layer at layer indices 3, 6, 9, 12. We notice that a
 557 single BST layer with state size $D = 16$ located closer to the middle of the whole Block Transformer
 558 stack, at index = 9, has the greatest effect on perplexity. This finding is inline with findings in prior
 559 work [36, 20].

560 In Table 3, we test if adding multiple BST layers yields improvements on performance. We start with
 561 BST layers with state size $D = 16$ at indices 0, 9. We follow by adding another BST layer at index 7
 562 for a total of three BST layers and then another at index 5, followed by another at index 12. Adding
 563 more BST layers lowers perplexity. However, the results seem to plateau at 5 BST layers. We note
 564 also that there is a 3.5% training step time increase for each added layer.

565 In Table 4 we train our models with different state sizes D . For the state size ablation, we use
 566 three BST layers at indices 0, 7, 9. We find that increasing D improves perplexity to the detriment of
 567 training speed (step time). For this reason, we chose $D = 16$ for Table 1 BST results.

568 D Limitations

569 While BST’s SSM layer allows the model to unroll and parallelize the recurrence that models long-
 570 term context between blocks of tokens, the SSM variants are reliant on efficient FFT operations. We
 571 have found that the FFT operation is an important speed bottleneck on TPUs that needs to be resolved
 572 to better scale BST to multiple layers and larger models. While we are still investigating the reasons,
 573 we found that JAX FFT was x4 faster on GPUs. Further, new SSM variants such as S5 [30] bypass
 574 FFT operations using a binary associative operator³. Our implementation is modular enough that we
 575 can simply plug in S5 or use other FFT implementations.

576 One of our assumption is that BST’s SSM layer is able to capture the right long-term dependency for
 577 each block. The SSM recurrence at step $T = t$ provides a summarized representation of previous
 578 steps for $T = 0$ to $T = t - 1$. However, a single vector representation may not be powerful enough
 579 to support all important long-term dependencies. Despite the perplexity improvements on long-range
 580 language modeling tasks, this assumption needs to be tested on other long range classification tasks
 581 such as Long Range Arena [32] as well. It is possible that our model can perform better if we feed to
 582 the attention layer W SSM representations that are chosen by a top-k retrieval operation, similar to
 583 the one in Memorizing Transformer [36].

584 E JAX Implementation of BST

585 Pseudocode 1 contains a function that implements convolution of multiple filters over the same input
 586 sequence using FFT and inverse FFT operations. Pseudocodes 2, 3 and 4 respectively implement
 587 context state collection of BST variants: Single-Head (SH), Multi-Head (MH) and Multi-Filter (MF).
 588 Finally, Pseudocode 5 runs the Block Transformer sublayer in parallel by feeding the context states
 589 to their corresponding block.

```
590 """Unstructured filters and convolutions."""
591
592 import jax
593 from jax import numpy as jnp
594 from einops import rearrange
595
596 win_length = 512 # (w)
597 seq_length = 4096 # (l)
598
599 def get_filters_unstruct(channels):
600     """Returns trainable filters and biases.
601
602     Args:
603         channels: number of filters.
604
605     Returns:
606         h: filter of shape (seq_length, channels, dim)
607         b: bias of shape (channels, dim)
608     """
609     t = jnp.linspace(0.0, 1.0, seq_length)
610     h = jnp.exp(- alpha * t) * dense(positional_emb(t))
611     b = get_bias()
612     return h, b
613
614 def multichannel_convolution(u, h, b):
615     """Multichannel convolution function.
616
617     Args:
618
```

³In JAX, this is equivalent to using `jax.lax.associative_scan`.

```

619         u: input of shape (seq_length, dim)
620         h: filters of shape (seq_length, channels, dim)
621         b: bias of shape (channels, dim)
622     """
623     h = rearrange(h, "l c d -> c d l")
624
625     fft_size = seq_length * 2
626     u_f = jnp.fft.rfft(x, n=fft_size)
627     h_f = jnp.fft.rfft(h, n=fft_size)
628
629     y = jnp.fft.irfft(h_f * x_f, n=fft_size, norm="forward")[
630         ..., :seq_length]          # (c, d, l)
631     y = y + x * b[..., None]       # (c, d, l)
632     y = rearrange(y, "c d l -> l d c")
633     return y
634

```

Pseudocode 1: Unstructured filters and convolutions.

```

635     """Context state collection for BST-SH variant."""
636
637     num_heads = 8          # (h)
638     num_states = 32       # (s)
639
640     # (SH): Single-Head
641     def SH_context_states(u):
642         """Single-Head Context Collection."""
643         h, b = get_filters_[unstruct/s4](channels=1)
644         y_1 = multichannel_convolution(u, h, b)
645         # y_1: (l, d, 1)
646
647         # lift to multiple heads
648         y_h = dense(y_1)
649         # y_h: (l, d, h)
650
651         context_states = jnp.split(
652             y_h, seq_length // win_length, axis=0)
653         return context_states # (l/w, w, d, h)
654

```

Pseudocode 2: Context state collection for BST-SH variants.

```

656     """Context state collection for BST-MH variant."""
657
658     # (MH): Multi-Head
659     def MH_context_states(u):
660         """Multi-Head Context Collection."""
661         h, b = get_filters_[unstruct/s4](channels=num_heads)
662         y_h = multichannel_convolution(u, h, b)
663         # y_h: (l, d, h)
664
665         context_states = jnp.split(
666             y_h, seq_length // win_length, axis=0)
667         return context_states # (l/w, w, d, h)
668

```

Pseudocode 3: Context state collection for BST-MH variants.

```

670     """Context state collection for BST-MF variant."""
671
672     # (MF): Multi-Filter
673     def MF_context_states(u):
674         """Multi-Filter Context Collection."""
675         h, b = get_filters_[unstruct/s4](channels=num_states)
676         y_s = multichannel_convolution(u, h, b)
677

```

```

678 # y_s: (l, d, s)
679 context_states = jnp.split(
680     y_s, seq_length // win_length, axis=0)
681 # context_states: (l/w, w, d, s)
682
683 # collect the last context states
684 context_states = context_states[:, -1, ...] # (l/w, d, s)
685 context_states = rearrange(
686     context_states, "lw d s -> lw s d")
687
688 # shift context states corresponding to windows
689 context_states = jnp.roll(context_states, 1, axis=1)
690
691 # replace the initial window with trainable weights
692 init_context = get_init_context(num_states) # (d, s)
693 context_states[0] = init_context
694
695 # lift to multiple heads
696 context_states = dense(context_states)
697
698 return context_states # (l/w, s, d, h)
699

```

Pseudocode 4: Context state collection for BST-MF variants.

```

700 """Block-State Transformer Layer."""
701
702 # Block Transformers are non-recurrent and parallelizable.
703 block_transformer = jax.vmap(BRecT.nonrecurrent_cell)
704
705 def BST(u):
706     """Block-State Transformer Layer."""
707     global MF # True if Multi-Filter, False otherwise (SH/MH)
708
709     # split inputs into windows (l/w, w, d)
710     u = jnp.split(u, seq_length // win_length, axis=0)
711
712     # collect context states from SSM outputs
713     context_states = [SH/MH/MF]_context_states(u)
714
715     # pass the contexts in place of recurrent states
716     y = block_transformer(
717         token_embeddings=u,
718         recurrent_state=context_states,
719         use_cross_attn_causal_mask=not MF,
720         use_cross_positional_emb=MF, # context IDs
721     )
722
723     return rearrange(y, "lw w d -> (lw w) d") # (l, d)
724

```

Pseudocode 5: Block-State Transformer Layer.