

Pangenome-Informed Language Models for Synthetic Genome Sequence Generation

Anonymous ACL submission

Abstract

Language Models (LM) have been extensively utilized for learning DNA sequence patterns and generating synthetic sequences. In this paper, we present a novel approach for the generation of synthetic DNA data using pangenomes in combination with LM. We introduce three innovative pangenome-based tokenization schemes that enhance long DNA sequence generation. Our experimental results demonstrate the superiority of pangenome-based tokenization over classical methods in generating high-utility synthetic DNA sequences, highlighting significant improvements in training efficiency and sequence quality.

1 Introduction

Public availability of genome datasets, such as the Human Genome Project (HGP) (Lander et al., 2001), the 1000 Genomes Project (Consortium et al., 2012), The Cancer Genome Atlas (TCGA) (Weinstein et al., 2013), GenBank (Benson et al., 2012), the International HapMap Project (Gibbs et al., 2003), the Human Pangenome Project (Liao et al., 2023), and the Telomere-to-Telomere project (Nurk et al., 2022), has been instrumental in advancing genomic research. However, large-scale genome sequencing remains costly and resource intensive due to the sophisticated equipment and computational resources required (Wetterstrand, 2021; Van Dijk et al., 2018).

Synthetic data generation offers a scalable alternative for genomic research. Specific tasks such as De Novo genome assembly (Tran et al., 2017, 2019; Yang et al., 2019) and genotype imputation (Brown and Browning, 2016) inherently involve the generation of unknown sequences, making them also suitable applications for synthetic data. A good generative model can significantly improve their accuracy and efficiency by predicting missing or incomplete segments.

Deep learning models are widely used in different tasks, even in processing genome sequences and related data (Yun et al., 2020; Kolesnikov et al., 2021; Kim and Kim, 2018; Elbashir et al., 2019). Although generative adversarial networks (GANs) have been explored for synthetic genome generation, their output is limited to short sequences (Bae et al., 2019; Gupta and Zou, 2018). LMs have shown their capability to generate synthetic natural languages that are almost indistinguishable from real data. The generated language text can be used to train other models (Kumar et al., 2020; Yoo et al., 2021; Hartvigsen et al., 2022), including those in the medical domain (Peng et al., 2023b; Guevara et al., 2024). Proven to be extraordinarily good at processing human language, LMs can also interpret and generate broader text, such as code for programming tasks (Chen et al., 2021), thereby pushing the boundaries of their application beyond strictly spoken language-based domains.

The Critical Challenge: DNA vs. NLP Tokenization Differences. While LMs present a promising alternative to understanding and generating long synthetic DNA sequences, effective tokenization of DNA sequences is crucial to leverage LMs. Applying LMs to DNA sequences faces a fundamental tokenization challenge that differs from Natural Language Processing (NLP). Unlike human language, which has natural word boundaries, semantic units, and grammatical structures, DNA is essentially a string with four letters (nucleotides: A, C, G, T) of billions of characters long and without inherent segmentation. Traditional NLP tokenization methods like Byte Pair Encoding (BPE) rely on frequency-based subword identification. Traditional DNA segmentation methods simply segment sequences into individual nucleotides or length k substrings. All these methods face a critical limitation as they lack awareness of the underlying biological structure that determines how genomic variations should be grouped and segmented.

Missing Opportunity: Structural Information from Pangenome Graphs. Pangenome graphs encode population-level variation patterns by comparing DNA sequences from a whole population and organizing DNA segments into nodes and connections that capture how genetic diversity manifests among individuals (detailed in §2.2). This graph structure naturally identifies biologically meaningful segmentation boundaries within the DNA sequences, offering a principled alternative to tokenizations that treat DNA as undifferentiated character strings. This opportunity has been largely overlooked in current DNA modeling approaches.

Therefore, we propose three novel pangenome graph-based tokenization schemes for LM-based synthetic data generation that leverage the structural information embedded in pangenome graphs to create biologically-informed segmentation.

This work presents the first comparative analysis of classical and pangenome-based tokenization schemes for LMs, specifically GPT-2 and Llama, in learning DNA sequence patterns and generating long synthetic sequences. Our findings reveal that the pangenome graph structure embeds significant information that enhances neural networks’ comprehension of DNA sequences, and can in cases reduce training time and improve scalability. Our contributions are as follows.

- **First**, we introduce three novel pangenome graph-based tokenization schemes that leverage genomic structure to provide biologically-informed segmentation, fundamentally different from frequency-based NLP tokenization approaches that ignore structural relationships.
- **Second**, we demonstrate through comprehensive experiments that our tokenization schemes significantly outperform three classical methods in training efficiency, predictive accuracy, and generation quality for both GPT-2 and Llama architectures, establishing clear computational and performance advantages.
- **Finally**, we establish the first systematic evaluation framework for pangenome-informed tokenization in synthetic DNA generation, providing evidence that structural graph information translates into improved biological utility through sequence alignment quality metrics.

The following paper is structured as follows: Section 2 covers background on synthetic genome generation, Section 3 details tokenization schemes, Section 4 outlines evaluation metrics, Section 5 presents experiments, Section 6 discusses related

work, and Section 7 concludes with limitations. A glossary is provided in Table 4 in §B.1.

2 Background

2.1 Language Models

Large language models are advanced artificial intelligence systems designed to understand and generate language text based on the data on which they have been trained. These models, such as Mistral (Jiang et al., 2023), Anthropic’s Claude (Anthropic, 2023), OpenAI’s GPT series (Radford et al., 2019; OpenAI, 2023), Google’s T5 (Raffel et al., 2020), Lamda (Thoppilan et al., 2022) and Gemini (Team et al., 2023), Meta’s OPT (Zhang et al., 2022), BLOOM (Le Scao et al., 2023) and Llama (Touvron et al., 2023a,b), etc., take advantage of vast amounts of textual information to learn patterns, nuances, and complexities of language. LMs can perform a variety of language-related tasks, including answering questions, translating languages, and even participating in casual conversations. Their ability to process and generate coherent and contextually appropriate responses makes them invaluable tools across multiple fields, from customer service and education to creative writing and technical support.

Model Choice with Computational Constraints. Due to computational limitations and the need for systematic tokenization comparison, we focus on accessible baseline models (90M parameter GPT-2 and Llama) trained from scratch rather than larger domain-specific DNA models. We selected these architectures for several key reasons: (1)*Tokenization control*: Tokenization flexibility to isolate the effects of different schemes by training from scratch; (2)*Generative capability*: We need autoregressive generation for synthetic sequence generation, whereas many smaller DNA-specific models (e.g., DNABERT-2 (Zhou et al., 2023), GENA (Fishman et al., 2023)) emphasize classifications; (3)*Feasibility*: Changing vocabularies alters the embedding matrix, requiring full retraining to isolate tokenization effects; DNA-specific models capable of generating long sequences exists (Nguyen et al., 2024a), but the billion-parameter structure makes it computationally too expensive for our academic research.

2.2 DNA Sequences and Pangenome Graphs

DNA basics. DNA carries genetic information and has four nucleotides: A, C, G, and T. A *genome*

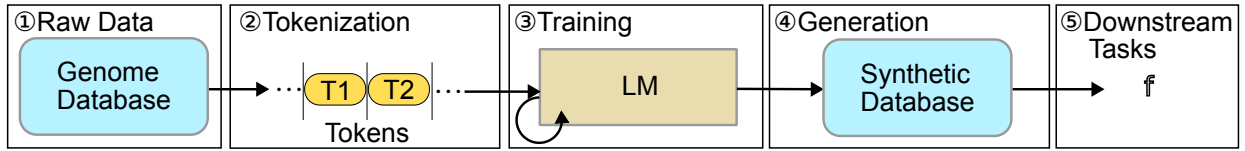


Figure 1: The whole pipeline of synthetic data generation and utilization.

is the complete DNA sequence of an organism, billions long in humans. Individuals of the same species share > 99% of this string, but small edits (single-base substitutions, insertions/deletions, structural rearrangements) occur throughout. Unlike natural language, DNA has no whitespace or punctuation; any segmentation is our modeling choice.

Pangenome graph. A “pangenome” aggregates the genomic content of many individuals of a species, where “pan” means “all”. A *pangenome graph* (Eizenga et al., 2020) encodes many genomes as a sequence graph, which is built by (whole-genome) multiple alignment: shared substrings in DNA are collapsed into **nodes**, substrings’ breakpoints become node boundaries, **edges** reconnect segments to preserve each genome’s order, and **paths** trace each individual (haplotype) through the graph. Figure 2 shows a toy graph of three sequences. Details of alignments are given in §4.2.

Why it is useful. The graph structure captures population-level variation while avoiding single-reference bias. Many downstream analyses operate on this structure, caring more about which branch an individual takes rather than on raw nucleotide counts. Importantly, this graph-based representation provides a more principled first-level segmentation that mimics natural language structure: just as words and phrases have meaningful boundaries in text, pangenome nodes represent biologically coherent segments where genetic variation naturally occurs, reflecting the underlying biological ‘grammar’. This natural segmentation provides an ideal foundation for tokenization. More details can be found in Appendix A.

2.3 Synthetic Genome Sequence Generation using LMs

We generate synthetic genomes with LMs through a five-step pipeline (See Figure 1): ① **Raw data** (§5.1) are collected; ② **Tokenization** (§3) converts sequences into model-ready tokens; ③ **LM training** fits a generative model with next-token prediction; ④ **Generation** (§5.1) samples synthetic sequences from the trained model; ⑤ **Downstream**

tasks (§4) assess the quality of the generation and how these sequences can be used.

3 Tokenization of a genome sequence

In this section, we describe the widely used tokenization schemes and our novel schemes. All tokenization schemes are illustrated in Figure 2.

3.1 Classical tokenizations

3.1.1 Genome-based Single Nucleotide Tokenization (GSNT)

GSNT is a straightforward method to tokenize genome sequences, previously applied in (Nguyen et al., 2024b; Schiff et al., 2024). Each nucleotide (A, C, G, T) is treated as an individual token. For example, the genome sequence “ACGTA” would be tokenized as “A”, “C”, “G”, “T”, and “A”.

3.1.2 Genome-based k -mer Tokenization (GKMT)

An alternative is GKMT, where k -mers, i.e. substrings of DNA of length k , are used as tokens. Depending on the stride, the k -mers may overlap or not overlap, and we focus on the non-overlapping k -mers in this work. Compared to GSNT, GKMT provides a longer effective context length, but is also highly sensitive to sequence mutations or errors: a single nucleotide insertion or deletion can change all subsequent tokens. For example, in Figure 2 GKMT, if the first A is missing, the segmentation will be “GCATGC TAGGCT...”, completely changing all tokens. A more detailed figure is shown in Figure 6 in §B.2.

3.1.3 Genome-based Byte Pair Encoding Tokenization (GBPET)

GBPET, also used in recent studies (Zhou et al., 2023), applies the BPE algorithm (Sennrich et al., 2016) to genome sequences. BPE begins with single nucleotide tokens and iteratively merges the most frequent pairs of adjacent tokens to create a vocabulary of longer subword-like tokens. However, BPE training requires too large computational resources if very long DNA sequences are given as inputs. Manual splitting of DNA text is needed in GBPET, which can harm the performance.

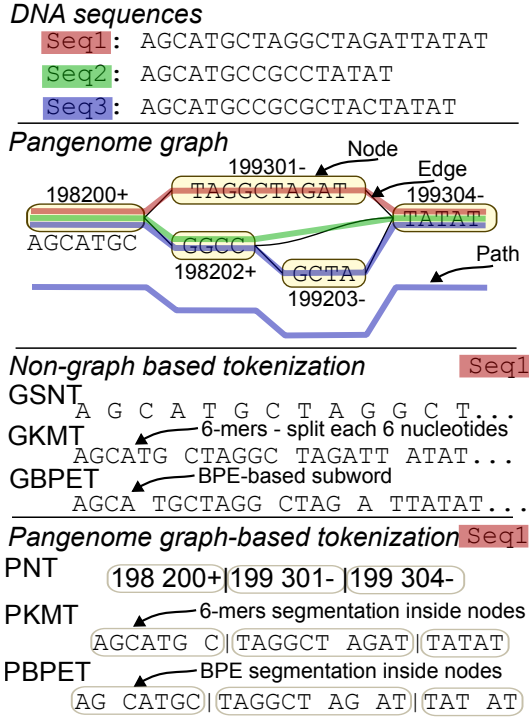


Figure 2: Three DNA sequences and a slice of a pangenome graph with nodes as bubbles (marked with ids), edges as lines. A valid path like the colored routes shows a single individual’s genome. Nodes can be interpreted as forward or reverse orientation, recorded this with +/- . The tokenization schemes output different segmented sequences of the red path.

3.2 Pangenome graph based tokenization

We introduce three novel pangenome graph-based tokenization schemes to overcome these limits:

3.2.1 Pangenome-based Node Tokenization (PNT)

Node IDs are assigned based on the pangenome graph construction process, which encodes both sequence content and its structural context to abstract numbers. PNT treats each graph node as a token. Because nodes encode both the DNA string and its position, identical strings at different locations receive different tokens, producing a very large vocabulary (~450k in our experiments vs. 50k in common NLP), which presents challenges for model training. We shrink the vocabulary size by splitting the node IDs into two parts (with an additional indicator for reversion, a common variation that causes the sequence to be replaced by its reverse complement). A drawback of PNT is its poor extensibility: adding new sequences requires rebuilding the graph and regenerating IDs, potentially changing the representations largely.

We therefore propose the next two schemes not

using IDs, and new sequences can be segmented using the original graph without a complete rebuild.

3.2.2 Pangenome-based k -mer Tokenization (PKMT)

PKMT first splits sequences at pangenome-graph node boundaries and then cuts each node string into non-overlapping k -mers as in GKMT (still $k=6$). Because the graph localizes insertions/deletions to specific nodes, PKMT is more tolerant to variations than GKMT. A main drawback compared with PNT is the loss of explicit positional/structural information from the graph, but using nucleotide strings rather than node IDs provides extensibility.

3.2.3 Pangenome-based BPE Tokenization (PBPET)

PBPET first splits sequences at pangenome-graph node boundaries as in PKMT, then runs Byte-Pair Encoding (BPE) that iteratively merges the most frequent adjacent symbol pairs and records the resulting merge rules (the “BPE merges”). We then apply the learned BPE merges to each node sequence, producing variable-length, high-frequency subwords while staying graph-aware via the node pre-segmentation. Unlike GBPET, this initial cut at node boundaries uses and preserves population-level structure from the graph.

4 Evaluating synthetic DNA generation quality

A main challenge in proving the utility of our schemes is evaluating the quality of the synthetic genome sequence generation. In our study, we use the prediction accuracy of the model to measure the quality of the generative model. Furthermore, we compare the similarity between synthetic and real genome sequences through sequence alignment.

4.1 Model Prediction Accuracy

- **Next token prediction accuracy:** measures how often the model correctly predicts the next token given the correct previous tokens, making it the primary metric for generative models. However, this does not fully reflect sequence accuracy when tokenization is not single nucleotide-based. Predicting “AAAAAC” or “GCTGCT” for the true k -mer token “AAAAAA” count both as simply incorrect.
- **Character-level prediction accuracy:** measures the percentage of nucleotides predicted correctly for each token. For example, predicting “AAAAAC” for the true token “AAAAAA” yields

an accuracy of 0.83, while predicting “GCTGCT” results in an accuracy of 0. Character-level accuracy measures how well each nucleotide is predicted instead of tokens, making it more consistent and fair across tokenization schemes.

4.2 Sequence Alignment Scores

The measurement of similarity between two genome sequences is done using sequence alignment, which is an essential process in many bioinformatic and computational biology tasks. Sequence alignment of DNA involves arranging the DNA sequences to identify regions of similarity. In our case, we use wfmash (Guarracino et al., 2025) where the wavefront algorithm (Marco-Sola et al., 2021) is primarily used for pairwise alignment between real and generated DNA sequences. Visualized results (introduced and shown in §5) and multiple scores can be used to evaluate the quality of the alignment.

An example of alignment between a reference sequence and a query sequence is shown in Figure 3.

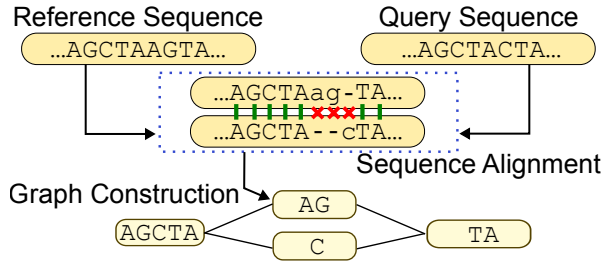


Figure 3: An alignment between two sequence and how it suggests graph nodes. Capitalized nucleotide and green links indicate matches; lowercase nucleotide and red crosses indicate no match; the dashes in the sequences represent the gaps during matching.

An alignment score of 0 indicates no similarity, while a score of 1 represents a perfect match. Alignment scores can be defined and computed in two primary ways (Figure 3 as an example):

- BLAST identity (BI): $7/10 = 0.7$. Defined as the number of matching bases in relation to the number of alignment columns.
- Gap-Compressed Identity (GI): $7/9 = 0.78$. Counting consecutive gaps as one difference.

Alignment scores reflect how well the aligned regions match, and we use **alignment percentage** to measure how much of a generated sequence aligns with references. Together, these metrics assess biological plausibility and generative quality.

Why alignment? We use alignment scores

as the primary evaluation metric for synthetic sequences, **rather than short-sequence classification tasks**, for three fundamental reasons:

(i) Dataset mismatch for classification tasks:

Most datasets of classification tasks, such as those used in DNABERT-2 (Zhou et al., 2023) and GENA (Fishman et al., 2023), evaluate the performance on short sequences (typically 10k or less), focusing on relatively more local characteristics. This length (with possible offset) is too short to build informative pangenome graphs that require long genomic regions of the same species.

(ii) Methodological mismatch: Classification tasks emphasize local motifs and short-range patterns, while pangenome graphs capture structural variants and long-range dependencies that are the core advantage of our approach.

(iii) Alignment as metric: Alignment scores quantify biological plausibility and consistently track downstream accuracy. Better alignment quality improves SNP/indel calling (Kosugi et al., 2013), and increases conserved gene detection, and aids comparative annotation (Sharma and Hiller, 2017). Alignment-based metrics better match practical genomic goals than generic divergence measures (Pillutla et al., 2021). Earlier discussions in computational biology also treated alignment score as a robust proxy for sequence similarity and usefulness (Frith, 2020; Durbin et al., 1998). In standard pipelines (e.g., reads → reference → variants), higher alignment of generated sequences to real genomes indicates they can stand in for real data (⑤ in §2.3; see Appendix E).

5 Experiments

5.1 Datasets and LM Choices

In our experiments, we used the human Major Histocompatibility Complex (MHC) region of chromosome 6 as our dataset, extracted from the PanGenome Graph Builder (PGGB) (Garrison et al., 2024) graph of Human Pangenome Reference Consortium (HPRC) year 1 assemblies (Liao et al., 2023). The MHC region was specifically chosen for its (1) high variation density with complex structural variants essential for pangenome graph construction; (2) long contiguous sequences necessary to capture the structural context that pangenome graphs encode, unlike typical classification datasets using short sequences (~ 10 kbp); and (3) population-level diversity across 126 samples with 447 million nucleotides total, providing

sufficient variation while maintaining biological realism. The dataset comprises 80% training samples and 20% test samples, with the reference genome temporarily used for hyperparameter tuning before final training.

We used 1024/2048 token context lengths for GPT-2/Llama respectively with Hugging Face 4.24.0 transformers library (Wolf, 2019). For BPE methods, we segment long genomes into 10k base pair (bp) sequences and set vocabulary size to 4096, following DNABERT2 (Zhou et al., 2023) specifications due to the computational limits of BPE tokenizer training on very long sequences. To mitigate memorization, we sample with sufficient randomness in generation (see Appendix C), and non-perfect alignment scores (<1) indicates that the model is not simply remembering and copying.

Table 1: Training time (hours) of each tokenization scheme on 90M models for 90 epochs.

Model	GSNT	GKMT	PKMT	GBPET	PBPET	PNT
GPT-2	56	11	15	17	24	7
LLaMA	20	5	6	9	12	3

5.2 Experiment Results

We trained the GPT-2 and Llama models on the dataset using four tokenization schemes: GSNT, GKMT, PNT and PKMT. Training was carried out for 90 epochs (§C.1 shows results with more epochs) with a batch size of 16/8 and 1024/2048-token sequences for GPT-2/Llama. The dataset comprises 124 DNA samples totaling 447 million nucleotides. Training times are shown in Table 1, obtained on 8 NVIDIA A5500 GPUs. Figure 4 displays token and character-level prediction accuracies. The final accuracies are shown in Table 2. PNT not included in the character-level accuracy figures due to the vague definition on predictions and targets with too varied lengths.

PNT demonstrated the fastest training time, while GSNT is generally the slowest due to its larger token set. BPE based method is slower than the k -mer based method but faster than GSNT. PNT reaches the best peak accuracy the fastest, while GKMT has the worst performance. GSNT initially trains much faster than PKMT for token prediction, but converges to a similar final accuracy. We will see how they perform differently in alignment. Despite having almost the same token tables, we can clearly tell that PKMT’s pangenome graph-aided segmentation helps the model to outperform the one trained by GSNT. The training of the PBPET

tokenizer takes around 20 seconds, while the training of GBPET tokenizer takes about 10 minutes, largely due to the larger sequence chunks, and they both have moderate training time.

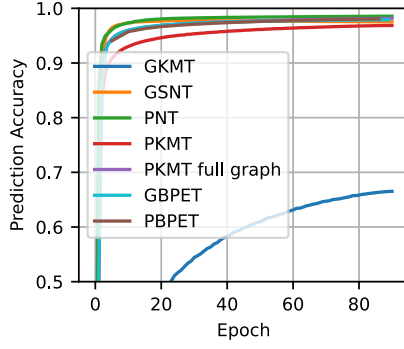
Table 2: Final accuracy of each tokenization scheme on 90M models trained for 90 epochs.

Model	GSNT	GKMT	PKMT	GBPET	PBPET	PNT
Token Prediction Accuracy						
GPT-2	97.1%	65.9%	96.9%	97.9%	98.0%	98.6%
LLaMA	98.7%	81.8%	97.7%	98.5%	98.6%	98.8%
Character-Level Accuracy						
GPT-2	97.1%	78.3%	97.9%	98.6%	99.0%	–
LLaMA	98.7%	85.3%	98.6%	99.0%	99.3%	–

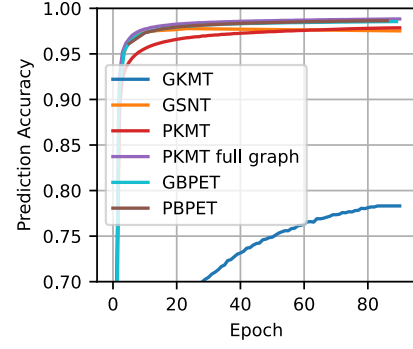
We report alignment results for GPT-2 generations across all tokenization schemes in Figure 5 (GKMT produces virtually no alignable sequence). The plot shows reference coordinates on the x -axis and individual generated sequences on the y -axis; each dot/segment marks a position in the generated sequence that aligns to the reference. After 90 epochs, only PNT yields long, contiguous alignments for GPT-2. Some sequences show no alignment at all, probably due to stochastic sampling for diversity and occasional misaligned patterns learned during training. LLaMA (see Appendix C) achieves similar token-level accuracy but exhibits sparser dots/dashes, indicating fewer matches overall. With PNT, LLaMA can initiate long runs, yet alignments tend to break off early, especially in high-variation regions (visible as dense dot clusters along the paths). With PKMT or PBPET, LLaMA cannot sustain long aligned sequences.

To quantify generation quality, we show the alignment scores of the generated sequences against the entire dataset (the best match of a query against the entire dataset) in Table 3, with the results for real data as a comparison. A variant-level check is given in §C.2. In addition to GI / BI scores, we show the alignment percentage, indicating the proportion of well-aligned sequences. The segment length refers to the size of the minimizer window during alignment. PNT achieves the highest alignment scores across all segment lengths, while GSNT performs the worst.

PNT demonstrates superior token-level prediction accuracy, while GKMT achieves the highest character-level accuracy in GPT-2 and closely rivals PNT in Llama. Traditional methods underperform, with GKMT achieving less than 70% accuracy and GSNT training significantly slower. The accuracy gap is more pronounced in alignment scores (Table 3), where PNT consistently excels



(a) Token prediction accuracy



(b) Character level prediction accuracy

Figure 4: Model prediction accuracies of all tokenization schemes during GPT-2 training.

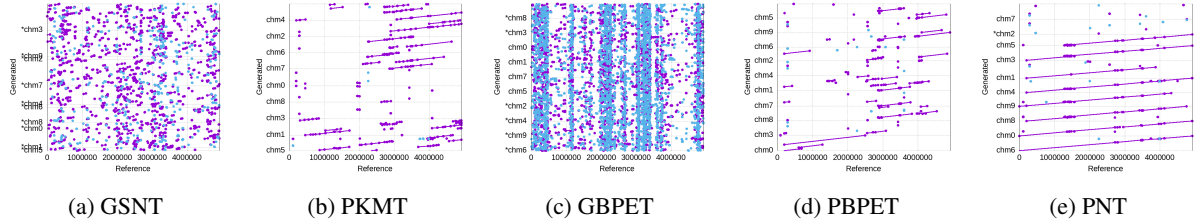


Figure 5: Alignment of a batch of GPT-2 generated sequences against the reference sequence. The X-axis represents the reference sequence; the Y-axis shows generated sequences. Longer lines indicate consistent alignment, and denser dots represent frequent short matches. Alignment results from Llama are presented in [Appendix C](#).

Table 3: Alignment percentages and weighted GI/BI scores of the 20 generated sequences per scheme for different segment lengths, aligned against the original dataset. Real data metrics are computed using 80% of samples as references and 20% as queries.

Segment	1k			20k			50k			200k		
GPT-2	Align %	GI	BI	Align %	GI	BI	Align %	GI	BI	Align %	GI	BI
GSNT	81.66	0.8712	0.9955	21.55	0.8834	0.9893	1.42	0.8323	0.9849	0.00	0.0000	0.0000
PKMT	52.96	0.9443	0.9856	50.34	0.9036	0.9932	47.87	0.8977	0.9936	8.82	0.8656	0.9919
GBPET	71.81	0.9873	0.9981	53.19	0.9105	0.9931	36.93	0.9041	0.9921	0.00	0.0000	0.0000
PBPET	44.75	0.9081	0.9914	42.03	0.9044	0.9943	42.29	0.9007	0.9935	9.39	0.9029	0.9955
PNT	89.34	0.9977	0.9999	31.01	0.9961	0.9990	33.27	0.9920	0.9985	36.96	0.9873	0.9982
LLaMA	Align %	GI	BI	Align %	GI	BI	Align %	GI	BI	Align %	GI	BI
GSNT	7.17	0.7927	0.9906	0.00	0.0000	0.0000	0.00	0.0000	0.0000	0.00	0.0000	0.0000
PKMT	34.45	0.9666	0.9960	21.19	0.8323	0.9907	6.85	0.8232	0.9876	0.00	0.0000	0.0000
GBPET	12.90	0.9543	0.9870	0.00	0.0000	0.0000	0.00	0.0000	0.0000	0.00	0.0000	0.0000
PBPET	33.00	0.9817	0.9969	6.41	0.8533	0.9861	5.86	0.8356	0.9878	0.00	0.0000	0.0000
PNT	28.80	0.9796	0.9958	5.97	0.9970	0.9984	8.49	0.9958	0.9987	15.02	0.9938	0.9977
Real data	99.97	0.9994	0.9999	69.23	0.9996	0.9999	61.37	0.9991	0.9997	50.67	0.9981	0.9993

with GI and BI scores of around 0.99 in segment lengths of 1k to 200k, closely mirroring the performance of real data. Although PKMT produces fewer high-quality sequences than GSNT that align with the reference, it achieves slightly higher alignment scores than GSNT in more settings and has a chance for relatively good generation for large segments. The newer non-pangenome-based method, GBPET, performs better in alignment score specifically under smaller segment length, but still lacks stable long-sequence generation compared with pangenome powered PBPET. PNT-generated sequences hold greater potential for applications re-

sembling real data, while others may require further refinement or model optimization. Llama overall shows the same trend, but lags behind GPT-2 in sequence generation, despite higher prediction accuracy and longer prompt length, likely due to greater performance degradation in very limited parameter numbers for continuous predictions. Llama specifically underperforms in non-pangenome based tokenization methods. Overall, PKMT performs better than GSNT (and GKMT), and PBPET performs better than GBPET, directly indicating the usefulness of involving pangenome graph structure in tokenization. One limitation we observe is that

pangenome-based models occasionally generate almost entire no match. Classical methods, although they generate fragmented pieces, do not completely miss. For PNT specifically, adding a small 20 token prompt will completely fix this issue.

Discussion. To our knowledge, this work is the first to compare the effectiveness of pangenome-based tokenization schemes with classical tokenization schemes for ML learning the pattern of DNA sequences; and also the one of the first to demonstrate the efficacy of LM in generating long DNA.

Our findings reveal that the pangenome graph structure embeds significant and meaningful information, improving neural networks’ understanding of DNA sequences. Our experiments demonstrate how this information can be effectively exploited by graph-based tokenization. The graph-aided segmentation of PKMT/PBPET provides more stable and learnable structural information compared to their classical counterpart, resulting in better overall generation quality. Our results underscore the trade-offs between computational cost and model performance, with pangenome graph-based tokenization schemes showing higher accuracy across tasks. Previous work (Liao et al., 2023) demonstrates how improved matching is the key point of the pangenome, which “aligns” with our use of the pangenome graph here.

6 Related work

In this section, we introduce two common genome tasks with the machine learning application. Table 8 in Appendix D summarizes this section.

6.1 Classification Tasks

Classification tasks are common in genomics, including (more details in Appendix E):

Variant Calling: ML models identify genetic variants such as SNPs and indels in genomes, linking them to diseases or traits. DeepVariant (Poplin et al., 2018), a CNN-based variant caller, outperforms traditional methods, influencing many others (Yun et al., 2020; Kolesnikov et al., 2021). Clairvoyante (Luo et al., 2019) excels in single-molecule sequencing (SMS), while Clair (Luo et al., 2020) offers faster RNN-based inference with fewer parameters, without sacrificing accuracy.

Gene Expression Analysis: ML models analyze gene expression data to reveal gene-disease relationships. Classical methods like KNN (Kim and Kim, 2018), linear/logistic regression (Han et al.,

2019), and SVMs (Wan et al., 2019) are used to predict driver genes or cancer risk. CNNs (Lyu and Haque, 2018; Elbashir et al., 2019) are also applied for cancer classification with RNA-seq data.

Beyond these, CNNs model protein binding (Alipanahi et al., 2015), cell type identification (Yao et al., 2019), and non-coding variants (Zhou and Troyanskaya, 2015). RNNs predict non-coding DNA functions (Quang and Xie, 2016) and RNA-protein binding preferences (Shen et al., 2020). Transformer models like DNA-BERT (Ji et al., 2021; Zhou et al., 2023; Dalla-Torre et al., 2023, 2025) provide strong contextual embeddings for molecular phenotype prediction but face context size limitations due to quadratic scaling. Recent models like Hyena (Nguyen et al., 2024b) and MambaDNA (Schiff et al., 2024) address these limitations with sub-quadratic scaling for longer contexts. More recent applications of DNA LM like MoDNA (An et al., 2022) for promoter prediction, and GENA (Fishman et al., 2023) for multiple tasks, both use traditional GKMT. Some papers like GPN-MSA (Benegas et al., 2024) for genome-wide variant effect prediction uses GSNT. DNABERT-2 (Zhou et al., 2023) and following work (Karollus et al., 2024) for evolutionary conservation and functional annotation prediction use BPE.

A recent paper (Zhang et al., 2024) presents a similar tokenization approach using pangenome graphs. Although both works independently develop this idea, ours differs by incorporating PNT and PBPET, and focusing on long-sequence generation. In contrast, their work handles shorter sequences (max 5000bp) with node-aided k -mer tokenization and focuses on classification tasks.

6.2 Generation Tasks

Synthetic Data Generation: Synthetic data mimics real data for privacy concerns. GANs have been used for synthetic medical data (Bae et al., 2019) and DNA sequences coding for proteins (Gupta and Zou, 2018), though limited by fixed output sizes. Some work (Avdeyev et al., 2023) utilizes transformers but with limited generation length, and a more recent large model (Nguyen et al., 2024a) shows generation of submillions in length with a certain level of genomic organization.

De Novo Genome Assembly: This involves reconstructing a genome from short DNA fragments without a reference. Deep learning has been applied to de novo peptide sequencing (Tran et al., 2017, 2019; Yang et al., 2019).

7 Limitations

PNT is the best performed scheme we proposed. However, it cannot tokenize sequences outside the nodes of the training pangenome graph, which makes it less general-purpose. We do have other two schemes that has no similar issue.

While our study focused on smaller models to establish a proof-of-concept for our tokenization scheme, we acknowledge that larger models may improve results but raise practical concerns around efficiency and resource use. Furthermore, emerging architectures designed for long-context processing (e.g., (Gu et al., 2021; Nguyen et al., 2024b,a; Gu and Dao, 2023; Peng et al., 2023a)) could potentially further enhance the performance of all tokenization schemes. These models, by enabling longer effective context windows, could improve both the understanding of long-range dependencies in DNA and the consistency of sequence generation. Although we believe that pangenome-based tokenization retains advantages in effective segmentation, such models may help close the performance gap for other tokenization methods. We agree that this is a valuable direction and suggest that future work explores scaling to larger models and incorporating long-context architectures to more fully assess their potential impact.

References

Babak Alipanahi, Andrew Delong, Matthew T Weirauch, and Brendan J Frey. 2015. Predicting the sequence specificities of dna-and rna-binding proteins by deep learning. *Nature biotechnology*, 33(8):831–838.

Weizhi An, Yuzhi Guo, Yatao Bian, Hehuan Ma, Jinyu Yang, Chunyuan Li, and Junzhou Huang. 2022. Modna: motif-oriented pre-training for dna language model. In *Proceedings of the 13th ACM international conference on bioinformatics, computational biology and health informatics*, pages 1–5.

Anthropic. 2023. [Claude 2](#). *Anthropic Blog*. Accessed: 2024-09-03.

Pavel Avdeyev, Chenlai Shi, Yuhao Tan, Kseniia Dudnyk, and Jian Zhou. 2023. Dirichlet diffusion score model for biological sequence generation. In *International Conference on Machine Learning*, pages 1276–1301. PMLR.

Ho Bae, Dahuin Jung, Hyun-Soo Choi, and Sungroh Yoon. 2019. Anomigan: Generative adversarial networks for anonymizing private medical data. In *Pacific Symposium on Biocomputing 2020*, pages 563–574. World Scientific.

Gonzalo Benegas, Carlos Albors, Alan J Aw, Chengzhong Ye, and Yun S Song. 2024. Gpn-msa: an alignment-based dna language model for genome-wide variant effect prediction. *bioRxiv*, pages 2023–10.

Dennis A Benson, Mark Cavanaugh, Karen Clark, Ilene Karsch-Mizrachi, David J Lipman, James Ostell, and Eric W Sayers. 2012. Genbank. *Nucleic acids research*, 41(D1):D36–D42.

Brian L Browning and Sharon R Browning. 2016. Genotype imputation with millions of reference samples. *The American Journal of Human Genetics*, 98(1):116–126.

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, and 1 others. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*.

1000 Genomes Project Consortium and 1 others. 2012. An integrated map of genetic variation from 1,092 human genomes. *Nature*, 491(7422):56.

Hugo Dalla-Torre, Liam Gonzalez, Javier Mendoza-Revilla, Nicolas Lopez Carranza, Adam Henryk Grzywaczewski, Francesco Oteri, Christian Dallago, Evan Trop, Bernardo P de Almeida, Hassan Sirelkhatim, and 1 others. 2023. The nucleotide transformer: Building and evaluating robust foundation models for human genomics. *bioRxiv*, pages 2023–01.

Hugo Dalla-Torre, Liam Gonzalez, Javier Mendoza-Revilla, Nicolas Lopez Carranza, Adam Henryk Grzywaczewski, Francesco Oteri, Christian Dallago, Evan Trop, Bernardo P de Almeida, Hassan Sirelkhatim, and 1 others. 2025. Nucleotide transformer: building and evaluating robust foundation models for human genomics. *Nature Methods*, 22(2):287–297.

Richard Durbin, Sean Eddy, Anders Størnøse Krogh, and Graeme Mitchison. 1998. Biological sequence analysis: Probabilistic models of proteins and nucleic acids.

Jordan M Eizenga, Adam M Novak, Jonas A Sibbesen, Simon Heumos, Ali Ghaffaari, Glenn Hickey, Xian Chang, Josiah D Seaman, Robin Rounthwaite, Jana Ebler, and 1 others. 2020. Pangenome graphs. *Annual review of genomics and human genetics*, 21:139–162.

Murtada K Elbashir, Mohamed Ezz, Mohanad Mohammed, and Said S Saloum. 2019. Lightweight convolutional neural network for breast cancer classification using rna-seq gene expression data. *IEEE Access*, 7:185338–185348.

Veniamin Fishman, Yuri Kuratov, Aleksei Shmelev, Maxim Petrov, Dmitry Penzar, Denis Shepelin, Nikolay Chekanov, Olga Kardymon, and Mikhail Burtsev. 2023. Gena-lm: a family of open-source foundational dna language models for long sequences. *bioRxiv*, pages 2023–06.

749	Martin C Frith. 2020. How sequence alignment scores	Alexander Karollus, Johannes Hingerl, Dennis Gankin,	804
750	correspond to probability models. <i>Bioinformatics</i> ,	Martin Grosshauser, Kristian Klemon, and Julien	805
751	36(2):408–415.	Gagneur. 2024. Species-aware dna language mod-	806
752	Erik Garrison, Andrea Guarracino, Simon Heumos,	els capture regulatory elements and their evolution.	807
753	Flavia Villani, Zhigui Bao, Lorenzo Tattini, Jörg	<i>Genome Biology</i> , 25(1):83.	808
754	Hagmann, Sebastian Vorbrugg, Santiago Marco-Sola,	Byung-Ju Kim and Sung-Hou Kim. 2018. Prediction of	809
755	Christian Kubica, and 1 others. 2024. Building	inherited genomic susceptibility to 20 common can-	810
756	pangenome graphs. <i>Nature Methods</i> , 21(11):2008–	cancer types by a supervised machine-learning method.	811
757	2012.	<i>Proceedings of the National Academy of Sciences</i> ,	812
758	Richard A Gibbs, John W Belmont, Paul Hardenbol,	115(6):1322–1327.	813
759	Thomas D Willis, Fuli L Yu, HM Yang, Lan-Yang	Alexey Kolesnikov, Sidharth Goel, Maria Nattestad,	814
760	Ch’ang, Wei Huang, Bin Liu, Yan Shen, and 1 others.	Taedong Yun, Gunjan Baid, Howard Yang, Cory Y	815
761	2003. The international hapmap project.	McLean, Pi-Chuan Chang, and Andrew Carroll. 2021.	816
762	Albert Gu and Tri Dao. 2023. Mamba: Linear-time	Deep trio: variant calling in families using deep learn-	817
763	sequence modeling with selective state spaces. <i>arXiv</i>	ing. <i>bioRxiv</i> , pages 2021–04.	818
764	<i>preprint arXiv:2312.00752</i> .	Shunichi Kosugi, Satoshi Natsume, Kentaro Yoshida,	819
765	Albert Gu, Karan Goel, and Christopher Ré. 2021. Effi-	Daniel MacLean, Liliana Cano, Sophien Kamoun,	820
766	ciently modeling long sequences with structured state	and Ryohei Terauchi. 2013. Coval: improving align-	821
767	spaces. <i>arXiv preprint arXiv:2111.00396</i> .	ment quality and variant calling accuracy for next-	822
768	Andrea Guarracino, Njagi Mwaniki, Santiago Marco-	generation sequencing data. <i>PloS one</i> , 8(10):e75402.	823
769	Sola, and Erik Garrison. 2025. wfmash: whole-	Varun Kumar, Ashutosh Choudhary, and Eunah Cho.	824
770	chromosome pairwise alignment using the hierarchi-	2020. Data augmentation using pre-trained trans-	825
771	cal wavefront algorithm .	former models. <i>arXiv preprint arXiv:2003.02245</i> .	826
772	Marco Guevara, Shan Chen, Spencer Thomas,	Eric S Lander, Lauren M Linton, Bruce Birren,	827
773	Tafadzwa L Chaunzwa, Idalid Franco, Benjamin H	Chad Nusbaum, Michael C Zody, Jennifer Baldwin,	828
774	Kann, Shalini Moningi, Jack M Qian, Madeleine	Keri Devon, Ken Dewar, Michael Doyle, William	829
775	Goldstein, Susan Harper, and 1 others. 2024. Large	Fitzhugh, and 1 others. 2001. Initial sequenc-	830
776	language models to identify social determinants of	ing and analysis of the human genome. <i>Nature</i> ,	831
777	health in electronic health records. <i>npj Digital</i>	409(6822):860–921.	832
778	<i>Medicine</i> , 7(1):6.	Teven Le Scao, Angela Fan, Christopher Akiki, El-	833
779	Anvita Gupta and James Zou. 2018. Feedback gan	lie Pavlick, Suzana Ilić, Daniel Hesslow, Roman	834
780	(fbgan) for dna: A novel feedback-loop architec-	Castagné, Alexandra Sasha Luccioni, François Yvon,	835
781	ture for optimizing protein functions. <i>arXiv preprint</i>	Matthias Gallé, and 1 others. 2023. Bloom: A 176b-	836
782	<i>arXiv:1804.01694</i> .	parameter open-access multilingual language model.	837
783	Yi Han, Juze Yang, Xinyi Qian, Wei-Chung Cheng,	Heng Li. 2018. Minimap2: pairwise alignment for	838
784	Shu-Hsuan Liu, Xing Hua, Liyuan Zhou, Yaning	nucleotide sequences. <i>Bioinformatics</i> , 34(18):3094–	839
785	Yang, Qingbiao Wu, Pengyuan Liu, and 1 others.	3100.	840
786	2019. Driverml: a machine learning algorithm for	Wen-Wei Liao, Mobin Asri, Jana Ebler, Daniel Doerr,	841
787	identifying driver genes in cancer sequencing studies.	Marina Haukness, Glenn Hickey, Shuangjia Lu, Ju-	842
788	<i>Nucleic acids research</i> , 47(8):e45–e45.	lian K Lucas, Jean Monlong, Haley J Abel, and 1	843
789	Thomas Hartvigsen, Saadia Gabriel, Hamid Palangi,	others. 2023. A draft human pangenome reference.	844
790	Maarten Sap, Dipankar Ray, and Ece Kamar. 2022.	<i>Nature</i> , 617(7960):312–324.	845
791	Toxigen: A large-scale machine-generated dataset for	Ruibang Luo, Fritz J Sedlazeck, Tak-Wah Lam, and	846
792	adversarial and implicit hate speech detection. <i>arXiv</i>	Michael C Schatz. 2019. A multi-task convolu-	847
793	<i>preprint arXiv:2203.09509</i> .	tional deep neural network for variant calling in sin-	848
794	Yanrong Ji, Zhihan Zhou, Han Liu, and Ramana V Davu-	gle molecule sequencing. <i>Nature communications</i> ,	849
795	luri. 2021. Dnabert: pre-trained bidirectional encoder	10(1):998.	850
796	representations from transformers model for dna-	Ruibang Luo, Chak-Lim Wong, Yat-Sing Wong, Chi-	851
797	language in genome. <i>Bioinformatics</i> , 37(15):2112–	Ian Tang, Chi-Man Liu, Chi-Ming Leung, and Tak-	852
798	2120.	Wah Lam. 2020. Exploring the limit of using a deep	853
799	Albert Q Jiang, Alexandre Sablayrolles, Arthur Men-	neural network on pileup data for germline variant	854
800	sch, Chris Bamford, Devendra Singh Chaplot, Diego	calling. <i>Nature Machine Intelligence</i> , 2(4):220–227.	855
801	de las Casas, Florian Bressand, Gianna Lengyel, Guil-	Boyu Lyu and Anamul Haque. 2018. Deep learning	856
802	laume Lample, Lucile Saulnier, and 1 others. 2023.	based tumor type classification using gene expression	857
803	Mistral 7b. <i>arXiv preprint arXiv:2310.06825</i> .	data. In <i>Proceedings of the 2018 ACM international</i>	858

859	<i>conference on bioinformatics, computational biology, and health informatics</i> , pages 89–96.	
860		
861	Santiago Marco-Sola, Juan Carlos Moure, Miquel Moreto, and Antonio Espinosa. 2021. Fast gap-affine pairwise alignment using the wavefront algorithm. <i>Bioinformatics</i> , 37(4):456–463.	
862		
863		
864		
865	Eric Nguyen, Michael Poli, Matthew G Durrant, Brian Kang, Dhruva Katrekar, David B Li, Liam J Bartie, Armin W Thomas, Samuel H King, Garyk Brixi, and 1 others. 2024a. Sequence modeling and design from molecular to genome scale with evo. <i>Science</i> , 386(6723):eado9336.	
866		
867		
868		
869		
870		
871	Eric Nguyen, Michael Poli, Marjan Faizi, Armin Thomas, Michael Wornow, Callum Birch-Sykes, Stefano Massaroli, Aman Patel, Clayton Rabideau, Yoshua Bengio, and 1 others. 2024b. Hyenadna: Long-range genomic sequence modeling at single nucleotide resolution. <i>Advances in neural information processing systems</i> , 36.	
872		
873		
874		
875		
876		
877		
878	Sergey Nurk, Sergey Koren, Arang Rhie, Mikko Rautiainen, Andrey V Bzikadze, Alla Mikheenko, Mitchell R Vollger, Nicolas Altemose, Lev Uralsky, Ariel Gershman, and 1 others. 2022. The complete sequence of a human genome. <i>Science</i> , 376(6588):44–53.	
879		
880		
881		
882		
883		
884	OpenAI. 2023. Gpt-4 technical report . <i>Preprint</i> , arXiv:2303.08774.	
885		
886	Donovan H Parks, Michael Imelfort, Connor T Skennerton, Philip Hugenholtz, and Gene W Tyson. 2015. Checkm: assessing the quality of microbial genomes recovered from isolates, single cells, and metagenomes. <i>Genome research</i> , 25(7):1043–1055.	
887		
888		
889		
890		
891	Bo Peng, Eric Alcaide, Quentin Anthony, Alon Albalak, Samuel Arcadinho, Stella Biderman, Huanqi Cao, Xin Cheng, Michael Chung, Matteo Grella, and 1 others. 2023a. Rwkv: Reinventing rnns for the transformer era. <i>arXiv preprint arXiv:2305.13048</i> .	
892		
893		
894		
895		
896	Cheng Peng, Xi Yang, Aokun Chen, Kaleb E Smith, Nima PourNejatian, Anthony B Costa, Cheryl Martin, Mona G Flores, Ying Zhang, Tanja Magoc, and 1 others. 2023b. A study of generative large language model for medical research and healthcare. <i>arXiv preprint arXiv:2305.13523</i> .	
897		
898		
899		
900		
901		
902	Krishna Pillutla, Swabha Swayamdipta, Rowan Zellers, John Thickstun, Sean Welleck, Yejin Choi, and Zaid Harchaoui. 2021. Mauve: Measuring the gap between neural text and human text using divergence frontiers. <i>Advances in Neural Information Processing Systems</i> , 34:4816–4828.	
903		
904		
905		
906		
907		
908	Ryan Poplin, Pi-Chuan Chang, David Alexander, Scott Schwartz, Thomas Colthurst, Alexander Ku, Dan Newburger, Jojo Dijamco, Nam Nguyen, Pegah T Afshar, and 1 others. 2018. A universal snp and small-indel variant caller using deep neural networks. <i>Nature biotechnology</i> , 36(10):983–987.	
909		
910		
911		
912		
913		
	Daniel Quang and Xiaohui Xie. 2016. Danq: a hybrid convolutional and recurrent deep neural network for quantifying the function of dna sequences. <i>Nucleic acids research</i> , 44(11):e107–e107.	914
		915
		916
		917
	Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.	918
		919
		920
	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>Journal of machine learning research</i> , 21(140):1–67.	921
		922
		923
		924
		925
		926
	Yair Schiff, Chia-Hsiang Kao, Aaron Gokaslan, Tri Dao, Albert Gu, and Volodymyr Kuleshov. 2024. Caduceus: Bi-directional equivariant long-range dna sequence modeling. <i>arXiv preprint arXiv:2403.03234</i> .	927
		928
		929
		930
	Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Neural machine translation of rare words with subword units. In <i>Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1715–1725.	931
		932
		933
		934
		935
	Virag Sharma and Michael Hiller. 2017. Increased alignment sensitivity improves the usage of genome alignments for comparative gene annotation. <i>Nucleic acids research</i> , 45(14):8369–8377.	936
		937
		938
		939
	Zhen Shen, Qinhu Zhang, Kyungsook Han, and De-Shuang Huang. 2020. A deep learning model for rna-protein binding preference prediction based on hierarchical lstm and attention network. <i>IEEE/ACM Transactions on Computational Biology and Bioinformatics</i> , 19(2):753–762.	940
		941
		942
		943
		944
		945
	Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, and 1 others. 2023. Gemini: a family of highly capable multimodal models. <i>arXiv preprint arXiv:2312.11805</i> .	946
		947
		948
		949
		950
		951
	Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, and 1 others. 2022. Lamda: Language models for dialog applications. <i>arXiv preprint arXiv:2201.08239</i> .	952
		953
		954
		955
		956
	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, and 1 others. 2023a. Llama: Open and efficient foundation language models. <i>arXiv preprint arXiv:2302.13971</i> .	957
		958
		959
		960
		961
		962
	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, and 1 others. 2023b. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> .	963
		964
		965
		966
		967
		968

969	Ngoc Hieu Tran, Rui Qiao, Lei Xin, Xin Chen, Chuyi	Susan Zhang, Stephen Roller, Naman Goyal, Mikel	1023
970	Liu, Xianglilan Zhang, Baozhen Shan, Ali Gh-	Artetxe, Moya Chen, Shuohui Chen, Christopher De-	1024
971	odsi, and Ming Li. 2019. Deep learning enables	wan, Mona Diab, Xian Li, Xi Victoria Lin, and 1	1025
972	de novo peptide sequencing from data-independent-	others. 2022. Opt: Open pre-trained transformer	1026
973	acquisition mass spectrometry. <i>Nature methods</i> ,	language models. <i>arXiv preprint arXiv:2205.01068</i> .	1027
974	16(1):63–66.		
975	Ngoc Hieu Tran, Xianglilan Zhang, Lei Xin, Baozhen	Xiang Zhang, Mingjie Yang, Xunhang Yin, Yining Qian,	1028
976	Shan, and Ming Li. 2017. De novo peptide sequenc-	and Fei Sun. 2024. Deepgene: An efficient founda-	1029
977	ing by deep learning. <i>Proceedings of the National</i>	tion model for genomics based on pan-genome graph	1030
978	<i>Academy of Sciences</i> , 114(31):8247–8252.	transformer. <i>bioRxiv</i> , pages 2024–04.	1031
979	Erwin L Van Dijk, Yan Jaszczyszyn, Delphine Naquin,	Jian Zhou and Olga G Troyanskaya. 2015. Predicting ef-	1032
980	and Claude Thermes. 2018. The third revolution	fects of noncoding variants with deep learning–based	1033
981	in sequencing technology. <i>Trends in Genetics</i> ,	sequence model. <i>Nature methods</i> , 12(10):931–934.	1034
982	34(9):666–681.		
983	Nathan Wan, David Weinberg, Tzu-Yu Liu, Katherine	Zhihan Zhou, Yanrong Ji, Weijian Li, Pratik Dutta, Ra-	1035
984	Niehaus, Eric A Ariazi, Daniel Delubac, Ajay Kan-	mana Davuluri, and Han Liu. 2023. Dnabert-2: Ef-	1036
985	nan, Brandon White, Mitch Bailey, Marvin Bertin,	ficient foundation model and benchmark for multi-	1037
986	and 1 others. 2019. Machine learning enables de-	species genome. <i>arXiv preprint arXiv:2306.15006</i> .	1038
987	tection of early-stage colorectal cancer by whole-		
988	genome sequencing of plasma cell-free dna. <i>BMC</i>	A Detailed Explanation of Pangenomes	1039
989	<i>cancer</i> , 19:1–10.	and Pangenome Graphs	1040
990	John N Weinstein, Eric A Collisson, Gordon B Mills,	A.1 Definition of Pangenome	1041
991	Kenna R Shaw, Brad A Ozenberger, Kyle Ellrott,	The term “pangenome” combines “pan” (meaning	1042
992	Ilya Shmulevich, Chris Sander, and Joshua M Stuart.	“all” or “every”) with the “genome,” referring to	1043
993	2013. The cancer genome atlas pan-cancer analysis	the complete collection of genetic material across	1044
994	project. <i>Nature genetics</i> , 45(10):1113–1120.	all individuals in a species or closely related group.	1045
995	K. A. Wetterstrand. 2021. Dna sequencing	Unlike traditional genomics, which relies on a sin-	1046
996	costs: Data from the nhgri genome sequenc-	gle representative genome sequence, pangenomics	1047
997	ing program (gsp). https://www.genome.gov/	embraces the full spectrum of genetic diversity.	1048
998	sequencingcostsdata . National Human Genome		
999	Research Institute.	A.2 How Pangenome Graphs are Built	1049
1000	T Wolf. 2019. Huggingface’s transformers: State-of-	The construction of a pangenome graph involves	1050
1001	the-art natural language processing. <i>arXiv preprint</i>	several computational steps:	1051
1002	<i>arXiv:1910.03771</i> .	Step 1: Sequence Collection Multiple high-	1052
1003	Hao Yang, Hao Chi, Wen-Feng Zeng, Wen-Jing Zhou,	quality genome assemblies are collected from di-	1053
1004	and Si-Min He. 2019. pnovo 3: precise de novo pep-	verse individuals within a species. These DNA	1054
1005	tide sequencing using a learning-to-rank framework.	sequences often undergo quality control to ensure	1055
1006	<i>Bioinformatics</i> , 35(14):i183–i190.	accuracy and completeness(e.g., ensure that the	1056
1007	Kai Yao, Nash D Rochman, and Sean X Sun. 2019. Cell	DNA sequences collected are from the same re-	1057
1008	type classification and unsupervised morphological	gions of different individuals).	1058
1009	phenotyping from low-resolution images using deep	Step 2: Multiple Sequence Alignment All col-	1059
1010	learning. <i>Scientific reports</i> , 9(1):13467.	lected sequences are aligned using sophisticated	1060
1011	Kang Min Yoo, Dongju Park, Jaewook Kang, Sang-	algorithms that can handle large-scale structural	1061
1012	Woo Lee, and Woomyeong Park. 2021. Gpt3mix:	variations. This alignment would compare each	1062
1013	Leveraging large-scale language models for text aug-	sequence with each other, and find the best over-	1063
1014	mentation. <i>arXiv preprint arXiv:2104.08826</i> .	all matching to determine the “shared” common	1064
1015	Taedong Yun, Helen Li, Pi-Chuan Chang, Michael F	subsequences and find variable regions where indi-	1065
1016	Lin, Andrew Carroll, and Cory Y McLean. 2020. Ac-	viduals differ. The overall goal of the alignment is	1066
1017	curate, scalable cohort variant calls using deepvariant	to establish a single, internally consistent multiple-	1067
1018	and glnexus. <i>Bioinformatics</i> , 36(24):5582–5589.	genome homology map. A successful alignment	1068
1019	Haoyang Zeng, Matthew D Edwards, Ge Liu, and	therefore maximises (i) the amount of sequence	1069
1020	David K Gifford. 2016. Convolutional neural net-	that can be confidently placed as shared columns	1070
1021	work architectures for predicting dna–protein binding.	and (ii) the biological plausibility of any gaps or	1071
1022	<i>Bioinformatics</i> , 32(12):i121–i127.		

rearrangements it introduces. In practice, this process is done by optimizing the alignment score as match rewards and mismatch penalties.

Step 3: Graph Construction

Once the multiple-sequence alignment is fixed, each homologous column (or contiguous run of columns) becomes a **node** containing that sequence, and **edges** connect nodes in the orders observed in the input genomes; where the graph branches, alternative **paths** capture the variant sequences. Sometimes the graph structure can then be optimized to minimize redundancy while preserving all genomic information.

A.3 Graph Components in Details

(1) **Nodes.** Each node contains:

- A DNA sequence (typically 100-10,000 base pairs)
- A unique identifier (node ID)
- Metadata about which genomes contain this sequence

Example: Node 123456 might contain the sequence “ATCGATCGAAGTC” and appear in 85% of the individuals in the pangenome. In the actual map they can be marked as either forward or reversed orientation.

(2) **Edges.** Edges represent adjacency relationships between nodes. Multiple incoming/outgoing edges indicates for alternative paths representing different variants

Example: Node A connects to both Node B and Node C, meaning some individuals have sequence A followed by B, while others have A followed by C.

(3) **Paths.** Each path represents one individual’s genome:

- A sequence of connected nodes
- Preserves the linear order of sequences in the original genome

A.4 Benefits of Pangenome Graphs

Pangenome graphs offer several advantages:

Comprehensive representation: All genetic variations in a population are captured in a single data structure. Further analysis results are less dependent on the choice of a single reference genome.

Context preservation: The graph maintains the genomic context around variations, which is crucial for understanding their functional impact.

Algorithmic efficiency: Many genomic analyses can be performed more efficiently in graph representation than on multiple individual genomes.

Scalability: New genomes can be incrementally extended to existing graphs without starting from scratch.

This comprehensive representation makes pangenome graphs particularly valuable for applications in domains such as population genomics or evolutionary biology, where understanding the full spectrum of genetic diversity is crucial.

B More on tokenization schemes

B.1 Glossary of Frequent Terms/Acronyms

Table 4: Glossary of Frequent Terms/Acronyms

Terms	Explanation
nucleotide	Units/Letters (A, G, C, T) of DNA
indel	Insertion and deletion
bp	Base pairs, A-T, C-G pair in DNA
GSNT	Genome-based Single Nucleotide Tokenization
GKMT	Genome-based k -mer Tokenization
GBPET	Genome-based BPE Tokenization
PNT	Pangenome-based Node Tokenization
PKMT	Pangenome-based k -mer Tokenization
PBPET	Pangenome-based BPE Tokenization

B.2 PKMT tokenization

We show how inserting and deletion can affect the GKMT in Figure 6.

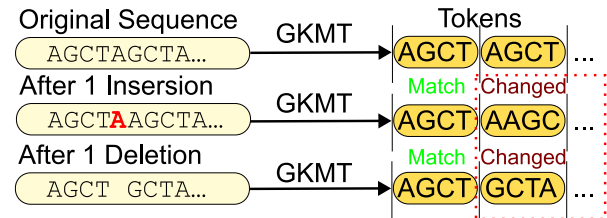
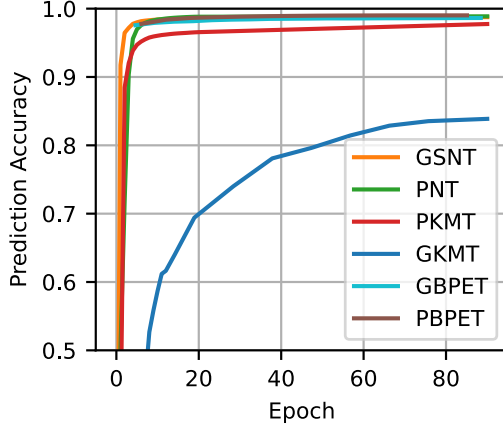


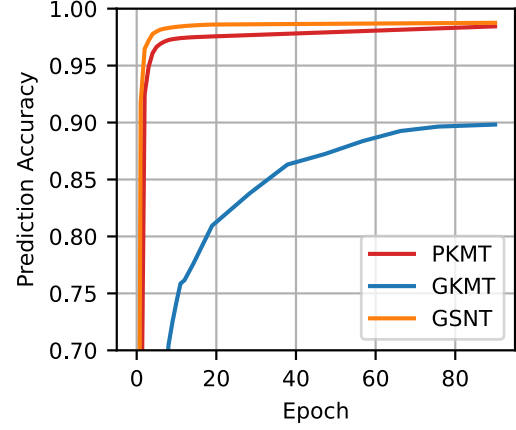
Figure 6: Insertion or Deletion of a single nucleotide change all following GKMT (stride equal to $k = 4$) tokens.

C More experiment details and results

The GPT-2 model uses the `gelu_new` activation function, consists of 12 transformer layers, each with 12 attention heads, and an embedding dimension of 768 with maximum prompt being 1024. The LLaMA model uses the `SiLU` activation function and consists of 6 transformer layers, each with 8 attention heads, and an embedding dimension of 768. It has an intermediate size of 4096, and supports sequences up to a maximum of 2048 positions. We used a grid search for the best hyperparameters. We use $3e-4$ (except $5e-4$ for GSNT - GPT-2 and $1e-4$ for LLaMA) learning



(a) Token prediction accuracy of the model across different training epochs



(b) Character-level prediction accuracy of the model across different training epochs

Figure 7: Model prediction accuracy of the four tokenization schemes during LLaMA training. PNT is excluded from the character-level accuracy plot due to the ambiguity in defining accuracy when predicted and target sequences differ in length.

Table 5: Alignment percentages and weighted GI/BI scores for segment lengths 5k and 100k.

Segment	5k			100k		
Model	Align %	GI	BI	Align %	GI	BI
GPT-2						
GSNT	59.49	0.8919	0.9910	0.00	0.0000	0.0000
PKMT	53.76	0.9015	0.9960	38.43	0.8928	0.9939
GBPET	63.20	0.9082	0.9961	14.62	0.9035	0.9884
PBPET	46.22	0.9019	0.9966	38.04	0.8927	0.9920
PNT	73.27	0.9970	0.9997	33.17	0.9945	0.9988
LLaMA						
GSNT	0.41	0.7414	0.9784	0.00	0.0000	0.0000
PKMT	30.82	0.8362	0.9917	0.00	0.0000	0.0000
GBPET	0.10	0.3070	0.9479	0.00	0.0000	0.0000
PBPET	24.75	0.8609	0.9916	0.00	0.0000	0.0000
PNT	25.79	0.9977	0.9997	10.96	0.9964	0.9990
Real data	97.97	0.9994	0.9999	60.44	0.9989	0.9997

rate, batch size 8/16 for GPT-2/Llama training; and $\text{topk}=10$, $\text{topp}=0.92$, $\text{topk_decend_min}=5$ for generation, which is also determined by grid search. Held-out accuracy is our proxy for interpolation: high token/character scores show generalization within the pangenome’s observed variation. Extrapolation is outweighed by interpolation due to species-level DNA homogeneity. We show additional alignment scores in Table 5 and the Llama accuracy in Figure 7. A clearer single query view of alignment is shown in Figure 8 for a single generated sequence, and the alignment figures for Llama are in Figure 10. Figure 9 shows a simple illustration of a small pangenome graph of the MHC data we use.

C.1 Effects of extensive training

During our experiment, we found that PNT, GBPET and PBPET did not benefit from more training epochs but GSNT and PKMT had the

potential for further improvement. We trained the better-performing GPT-2 model on half of the training dataset for an extra 200 epochs, keeping other parameters the same to further investigate the best possible performance these two tokenization schemes can provide. The token prediction accuracy increased by about 0.4% for PKMT and 0.3% for GSNT, which is marginal, but we observed significant improvements in generation quality for both methods in alignment score. While prediction accuracy gains may appear small, they have a compounding effect during generation, where errors accumulate across long sequences. Accuracy reflects only top-1 correctness for the next token, whereas generation samples probabilistically from the top candidates, making it more sensitive to distributional improvements. The results are shown in Figure 11 and Table 6. Both methods achieved slightly higher alignment scores and aligned length, especially with larger segments. Both tokenization

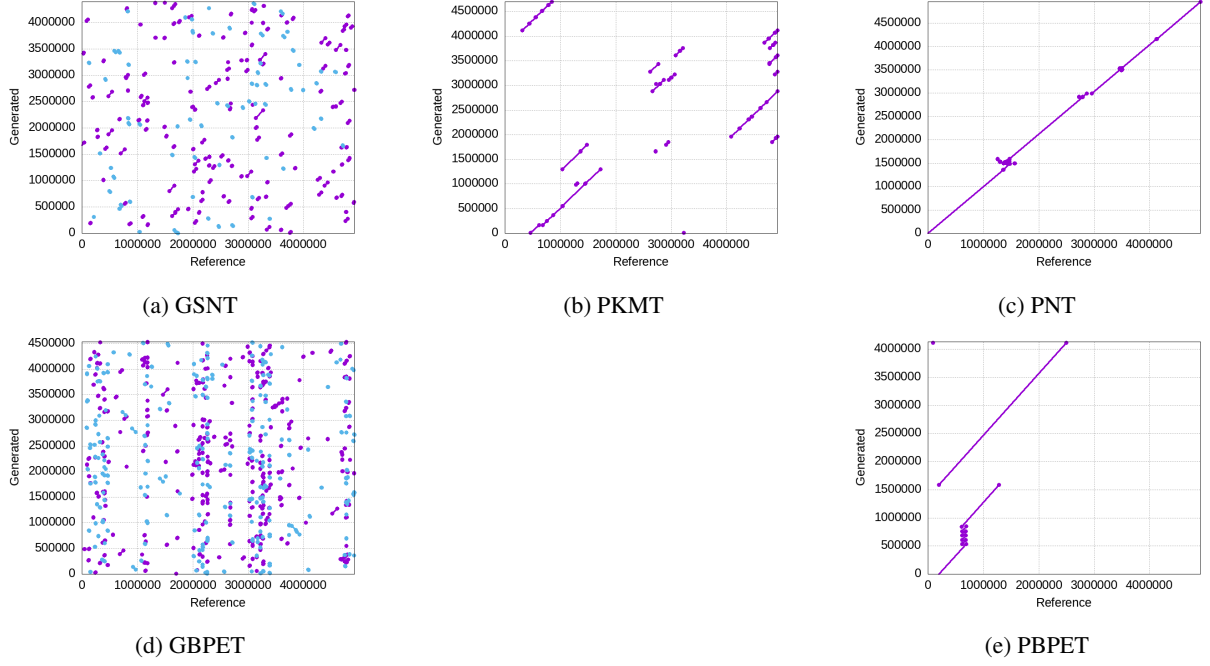


Figure 8: Alignment of a single generated sequence against the reference. Longer lines represent continuous alignment regions, while scattered dots show shorter matching fragments.

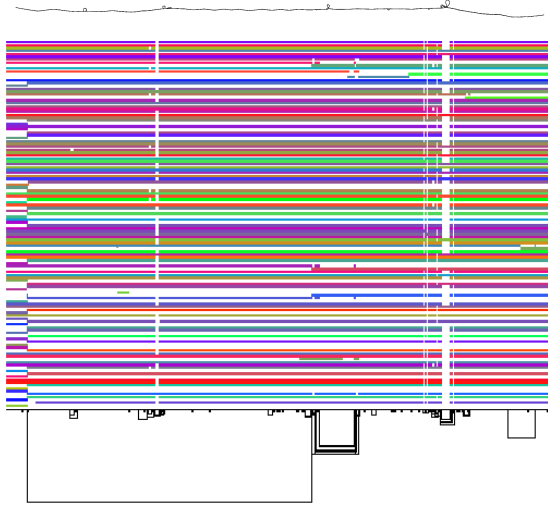


Figure 9: The pangenome graph of the human MHC region of chromosome 6 of the PGGB graph of HPRC year 1 assemblies, with 2D graph visualization (above) and matrix view (below). The circled in the 2D graph and the gaps in the matrix view indicate mutations.

schemes still underperformed compared to PNT, even after extensive training. Figure 11 additionally clearly shows that GKMT generates relatively longer sequences with more longer lines.

The higher utility of the extensively trained model indicates that substantial investment in computational power has its potential.

C.2 Downstream task: Variant calling

Variant calling evaluates whether the model captures *heterogeneity* (population-level variation) rather than merely memorizing conserved, homogeneous regions. More details are described in §E.1; here we summarize the experiment setup.

For each tokenization scheme, we sampled 20 sequences generated by the GPT-2 model (chosen because it achieved the best alignment-based performance). We then combined these synthetic sequences with the single training reference path and built a pangenome graph using PGGB (Garrison et al., 2024). Variants were called for each synthetic sequence against the reference path and compared to the “truth” set obtained by calling variants for all real (non-reference) genomes in the original pangenome graph.

Let T be the set of variants from real genomes (ground truth) and G the set from generated sequences. After standard normalization (left-normalization, multiallelic splitting, etc.), a *match* is defined when a variant in G is equivalent to one in T (e.g., same chromosome, position, and REF/ALT, or an accepted equivalence rule).

- **True positives (TP):** variants in G that match a variant in T .
- **False positives (FP):** variants in G with no match in T .
- **False negatives (FN):** variants in T with no match in G .

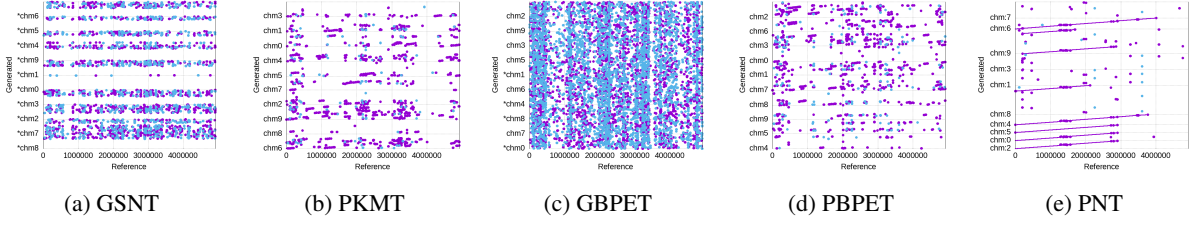


Figure 10: Alignment of a batch of LLaMA-generated sequences against the reference. The X-axis is the reference, and the Y-axis shows the generated sequences. Longer lines indicate consistent alignment, and denser dots indicate frequent short matches.

Table 6: Alignment percentages and weighted GI/BI scores of the 20 generated sequences each scheme for different segment lengths of the generated sequences with extensively trained GPT-2 model, against the test set as reference.

Segment	1k			5k			20k		
	Align %	GI	BI	Align %	GI	BI	Align %	GI	BI
GSNT	90.67	0.8818	0.9972	75.52	0.8922	0.9926	42.17	0.8916	0.9920
PKMT	81.42	0.9842	0.9978	81.87	0.9027	0.9969	79.74	0.9044	0.9956
Segment	50k			100k			200k		
	Align %	GI	BI	Align %	GI	BI	Align %	GI	BI
GSNT	9.76	0.8801	0.9916	0.00	0.0000	0.0000	0.00	0.0000	0.0000
PKMT	72.52	0.9011	0.9940	65.51	0.8936	0.9943	16.19	0.8883	0.9935

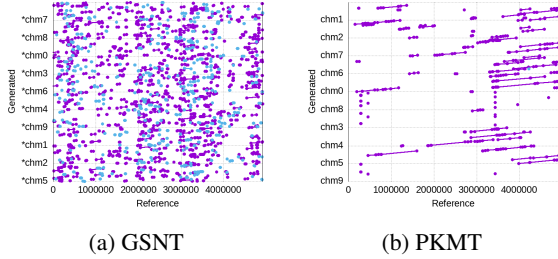


Figure 11: Alignment of a batch of generated sequences (after extensive GPT-2 training) against training sequences. The X-axis is the reference; the Y-axis contains generated sequences. Longer lines indicate consistent alignments, while denser dots reflect frequent short matches.

outperform classical baselines in all metrics, indicating they help the model generate biologically plausible variants (true positives). They more efficiently capture population diversity rather than just conserved sequence.

Method	Precision (%)	Recall (%)	F1-score (%)
GSNT	29.5	16.4	21.0
PKMT	65.7	28.2	39.4
GBPET	58.1	22.9	32.8
PBPET	74.6	29.7	42.4
PNT	69.3	23.3	34.8

Table 7: Variant-calling performance (precision, recall, and F_1) for each tokenization method.

Precision (of what we predicted, how much was correct?):

$$\text{Precision} = \frac{|TP|}{|TP| + |FP|}.$$

Recall (of what was true, how much did we recover?):

$$\text{Recall} = \frac{|TP|}{|TP| + |FN|}.$$

The F_1 score balances both:

$$F_1 = \frac{2 \cdot \text{Prec} \cdot \text{Rec}}{\text{Prec} + \text{Rec}} = \frac{2|TP|}{2|TP| + |FP| + |FN|}.$$

All scores are reported as percentages in Table 7. Recall values are modest because only 20 synthetic sequences are contrasted against a truth set derived from 126 real genomes; consequently, many true variants simply never appear in the generated subset. Nonetheless, all pangenome-based tokenizers

D Summarizing related work

Here we provide a table to summarize our discussion in §6, with a detailed list of the related work of ML/DL doing genomic tasks.

E Alignment scores and downstream tasks

Alignment-based evaluations provide a more direct assessment of how well synthetic data supports real-world genomic applications. For example, datasets like those from the Human Pangenome Project depend heavily on alignment-based metrics to assess data quality and interpret genetic variation. Read alignment to a reference genome followed by variant calling is a widely adopted pipeline, and here alignment consistency and accuracy are critical. In this context, alignment scores are not only practical but also well-recognized within the genomics community as meaningful indicators of quality.

In this section, we introduce two essential tasks to show how alignment scores can determine the utility of sequences, and how synthetic sequences can play a role.

E.1 Variant calling

Read alignment and variant calling are foundational tasks in bioinformatics pipelines, especially in genome resequencing studies. In this process, DNA reads generated by sequencing technologies are aligned to a reference genome to reconstruct the original genetic material and identify variants (e.g., calling the inserting and deletion in the bottom two sequences when compared with the top reference in Figure 6). Determining an accurate alignment is critical because downstream variant calling algorithms rely on these mappings to compare the sample DNA against the reference. Numerous tools have been developed to perform this task efficiently and accurately, including Minimap2 (Li, 2018) and wfmash (Guarracino et al., 2025). Most work in §6.1 measure the alignment in their experiment.

A high alignment score indicates a strong match between the sequenced read and a region in the reference genome, minimizing mismatches, gaps, or ambiguous placements. This is essential to identify true variants confidently, ruling out sequencing errors or misalignments. An incorrect alignment may map a query DNA sequence to the wrong location in the reference genome, leading to wrong variant calls. An example is given in Figure 12. Synthetic sequences can serve as references in variant calls or generate potential variant combinations that are not observed in natural samples.

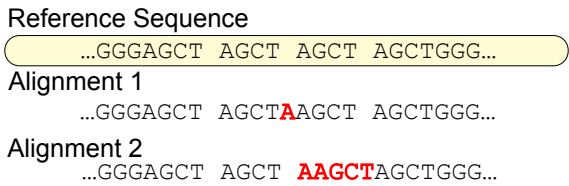


Figure 12: Two possible alignment of a sequence to a reference sequence. Alignment 1 calls for one insertion while Alignment 2 calls for 4 deletion then 5 insertion. Alignment 1 will have higher alignment scores with more matched nucleotides, and is considered a better alignment. Therefore the variant calling based on Alignment 1 is considered better than Alignment 2.

E.2 De novo assembly

De novo assembly reconstructs a genome from short sequencing reads without relying on a ref-

erence genome. This process stitches overlapping reads into contiguous sequences (contigs) or scaffolds, aiming to rebuild the original genome as accurately as possible. Since there is no reference during assembly, evaluation is typically performed by aligning the assembled contigs back to a trusted reference genome, or comparing them to known markers or conserved genes.

A high alignment score here indicates that the assembler has likely reconstructed a biologically accurate sequence. This suggests high contiguity, low error rates, and minimal misassemblies. Low alignment scores often signal fragmented or misassembled regions. Synthetic sequences can act as trusted reference, improving the assembly.

Many utility metrics used in existing genome modeling studies are fundamentally rooted in sequence alignment. For example, in recent work such as (Nguyen et al., 2024a), tools like CheckM (Parks et al., 2015) are used to report quality metrics, including gene density and stop codon frequencies. These tools rely on foundational components like profile Hidden Markov Models (pHMMs) that are directly constructed from multiple sequence alignments, with alignment quality and consistency playing a central role in shaping their parameters and performance. In this context, a high alignment score indicates strong homology or functional similarity between the generated sequence and known sequences, providing evidence of biological plausibility.

Table 8: DL models used in genome tasks.

Job Type	Paper	Task	Architecture	Input
Classification	(Poplin et al., 2018; Yun et al., 2020; Kolesnikov et al., 2021)	Variant Calling	CNN	hundreds of base pairs
	(Luo et al., 2019)	Variant Calling	CNN	hundreds of base pairs
	(Lyu and Haque, 2018; Elbashir et al., 2019)	Cancer Prediction	CNN	RNA-seq gene expression data
	(Alipanahi et al., 2015)	Protein Binding	CNN	10-100 nucleotides & binding specificities
	(Zeng et al., 2016)	Protein Binding	CNN	10-100 base pairs & binding specificities
	(Yao et al., 2019)	Cell Type Identification	CNN	cell images
	(Zhou and Troyanskaya, 2015)	Non-coding DNA function prediction	CNN	1k base pairs
	(Luo et al., 2020)	Variant Calling	RNN	binary alignment map (BAM)
	(Shen et al., 2020)	RNA-protein binding preference	LSTM	embedded k -mers
	(Quang and Xie, 2016)	Non-coding DNA function prediction	CNN/BLSTM	one hot encoded nucleotides
	(Kim and Kim, 2018)	Cancer Prediction	KNN	SNP genotype syntaxes (8-mers)
	(Han et al., 2019)	Cancer Prediction	Rao score	Mutation Annotation Format (MAF)
	(Wan et al., 2019)	Cancer Prediction	SVM	Human EDTA plasma samples
	(Ji et al., 2021; Zhou et al., 2023)	Molecular Phenotype Prediction	Transformer	tokenized k -mers
	(Dalla-Torre et al., 2023)	Molecular Phenotype Prediction	Transformer	tokenized k -mers
	(Nguyen et al., 2024b)	5-way Species Classification	Transformer	single nucleotide tokens
	(Schiff et al., 2024)	Genome Tasks	Mamba	single nucleotide tokens
	(Luo et al., 2019)	Variant Calling	CNN	Hundreds of base pairs
	(An et al., 2022)	Promoter Prediction	Transformer	6-mers of up to 512bp
	(Karollus et al., 2024)	Evolutionary Conservation / Functional Annotations	Transformer	6-mers for 128bp sequences
	(Fishman et al., 2023)	Multiple Tasks	Transformer	BPE tokens, up to 36000bp sequence
	(Benegas et al., 2024)	Genome-wide Variant Effect Prediction	Transformer	GSNT for 128bp sequences
	(Dalla-Torre et al., 2025)	Multiple Prediction Tasks	Transformer	Thousands of k -mer tokens
Generation	(Tran et al., 2017)	De novo peptide sequencing	LSTM/CNN	tandem mass spectrometry (MS/MS) Spectrum
	(Tran et al., 2019)	De novo peptide sequencing	LSTM/CNN	data-independent acquisition (DIA) mass spectrometry data
	(Yang et al., 2019)	De novo peptide sequencing	learning-to-rank	tandem mass spectrometry data
	(Bae et al., 2019)	Synthetic Medical Data	GAN	medical data
	(Gupta and Zou, 2018)	Synthetic DNA Sequences	GAN	DNA sequences
	(Avdeyev et al., 2023)	Synthetic DNA Sequences	Transformer	Up to 1024 base-pairs
	(Nguyen et al., 2024a)	Synthetic DNA Sequences	Transformer	Up to 131072 base-pairs