

ZSMILES: an approach for efficient SMILES storage for random access in Virtual Screening

Gianmarco Accordi*, Davide Gadioli*, Giorgio Seguini*, Andrea R. Beccari†, and Gianluca Palermo*

*Dipartimento di Elettronica Informatica e Bioingegneria - Politecnico di Milano, Milano, Italy

†EXSCALATE - Dompé Farmaceutici SpA, Naples, Italy

Email: *{firstname.lastname}@polimi.it †andrea.beccari@dompe.com

Abstract—Virtual screening is a technique used in drug discovery to select the most promising molecules to test in a lab. To perform virtual screening, we need a large set of molecules as input, and storing these molecules can become an issue. In fact, extreme-scale high-throughput virtual screening applications require a big dataset of input molecules and produce an even bigger dataset as output. These molecules’ databases occupy tens of TB of storage space, and domain experts frequently sample a small portion of this data. In this context, SMILES is a popular data format for storing large sets of molecules since it requires significantly less space to represent molecules than other formats (e.g., MOL2, SDF).

This paper proposes an efficient dictionary-based approach to compress SMILES-based datasets. This approach takes advantage of domain knowledge to provide a readable output with separable SMILES, enabling random access. We examine the benefits of storing these datasets using ZSMILES to reduce the cold storage footprint in HPC systems. The main contributions concern a custom dictionary-based approach and a data preprocessing step.

From experimental results, we can notice how ZSMILES leverage domain knowledge to compress $\times 1.13$ more than state of the art in similar scenarios and up to 0.29 compression ratio. We tested a CUDA version of ZSMILES targeting NVIDIA’s GPUs, showing a potential speedup of $7\times$.

I. INTRODUCTION

The drug discovery process is a lengthy and costly process for a pharmaceutical company [1]. The process produces molecules (called *ligands*), which are more likely to interact with at least one binding site (called *pocket*) on a protein (called the screening’s *target*) [2]. Ideally, these interactions inhibit the target’s activity, leading to a favorable therapeutic effect. The process comprises in-silico, in-vitro, and in-vivo stages.

Virtual screening is an in-silico stage, which can improve the success rate [3] of the drug discovery pipeline by selecting chemical compounds with a higher likelihood of interacting with the target protein. We screen compounds from a large dataset to ensure the most likely interactions [4], [5]. Virtual screening eliminates the molecules with low binding affinity for in-vitro and later in-vivo testing [6].

In a virtual screening campaign, evaluating the initial compounds’ interaction strength against multiple target proteins is common. This cartesian product increases the required computation effort. Moreover, the evaluation of each ligand-protein pair is independent of the others, making the problem

embarrassingly parallel. For these reasons, supercomputers are the ideal target for extreme-scale virtual screening campaigns.

One challenge of virtual screening campaigns is the storage requirement [7], [8], to describe the input chemical library and the output, which usually decorates the input with the strength of their interactions. For example, the screening data of a virtual screening campaign on CINECA’s Marconi100 [8] was approximately 72 TB. Not all screening data is accessed on a daily basis, but rather, domain experts sample this chemical space to create a smaller subset. To mitigate this problem, it is common to encode molecules using the SMILES format, which describes a molecule using a single line of ASCII characters. Given the rising popularity of large virtual screening campaigns and SMILES format, storing them efficiently is of general interest.

This work presents ZSMILES, a methodology to reduce the SMILES storage footprint of extreme-scale virtual screening applications. The proposed approach employs a dictionary-based compression (and decompression). Leveraging domain knowledge to reduce SMILES’s storage footprint is one of the contributions of this paper. Our approach takes advantage of SMILES format specifications to get better compression. In particular, we have addressed the dictionary generation and data preparation phase. Moreover, we have other constraints: domain experts have to access these databases easily without, for example, the burden of handling binary characters. Domain experts must also cut and combine SMILES databases: we defined the shared dictionary to be input-independent and SMILES to remain separable. Thus, our approach employs a single fixed dictionary to compress any set of SMILES and enhance maintainability and compatibility. SMILES separability implies each SMILES to be placed on different lines. ZSMILES’s output still provided each SMILES on the same input line number. In the following with random access, we refer to maintaining input SMILES order in output, thus for SMILES to remain separable. We implemented ZSMILES in a serial C++ version targeting CPUs. Large virtual screening experiments are usually accelerated using GPUs. Thus, we have tested a parallel CUDA version targeting GPUs.

We evaluate the impact of domain-specific optimizations on performance. ZSMILES compression ratios are evaluated with different optimizations enabled and dictionaries trained on different datasets of SMILES. We also analyze ZSMILES’s



Fig. 1: Graphical representation of Vanillin on the left, while on the right, its SMILES representation.

implementations execution times.

The remainder of this article is structured as follows: initially, we report the background Section II. Subsequently, we provide a state-of-the-art analysis in Section III. Next, Section IV outlines the proposed methodology and illustrates the accelerated approach. Finally, an analysis of the results acquired in section Section V is presented, with the conclusions given in Section VI. Moreover, abbreviations and technical terms are explicitly defined when first used to facilitate clear communication.

II. BACKGROUND

The methodology proposed in this study focuses on SMILES strings, which can describe molecules using a single-line notation that encodes only its topological information using UTF-8 ASCII characters [9]. Therefore, the proposed approach falls in the category of short-string compression algorithms. This format provides an excellent balance between human-readable and computationally parsable notations. The main idea is to start the molecule description from a terminal heavy atom and write all the other heavy atoms attached to the starting one in sequence. It uses round brackets to handle branches in the molecule structure. When we encounter a ring, we remove an edge and assign a numerical ID to the attached atoms, and then we resume the encoding procedure. Therefore, this numerical ID can identify all the molecule rings. The SMILES format has additional rules to encode chemical information, such as upper-case chemical symbols for non-aromatic atoms and lower-case chemical symbols for aromatic ones. From a SMILES string, it is possible to add hydrogens and compute the 3D displacement of their atoms at runtime.

For example, Figure 1 shows a graphical representation of the Vanillin, COc1cc(C=O)cc1O in SMILES format. In this example, the SMILES representation starts from the bottom part of the molecule, i.e., from the lowermost carbon. Thus, our SMILES notation starts with the *C* (since this carbon is not aromatic). Following the molecule structure, the next heavy atom is a non-aromatic oxygen, expanding the SMILES string to *CO*. The next atom is an aromatic carbon that belongs to a ring. To represent it, we remove the edge on its left (represented by the dashed lines), assign ID 1, and add the carbon to SMILES notation *COc1*. Then, we resume the encoding structure, including two aromatic carbons, i.e.,

COc1cc. At this point, we reach a branch that we need to handle using round brackets. For simplicity, we assume to encode the right branch first, so our SMILES representation becomes *COc1cc(C=O)* (the = symbol stands for a double covalent bond). Then, we also encode the left branch, which includes all the remaining heavy atoms, i.e., *COc1cc(C=O)cc1O*.

Please notice how the SMILES description is not unique; it depends on how we topologically order the molecule structure. Moreover, the ring enumeration does not have to be unique. The only requirement is that nested ring descriptions should not have overlapping IDs to prevent ambiguities in their representation.

III. STATE OF THE ART

Dictionary-based compression algorithms [10] work by parsing the input while attempting to match substrings against a predefined set within a dictionary. When a match occurs, the algorithm substitutes the substring in the input with the correspondent symbol of the dictionary entry: if the latter is shorter than the former, we have reduced the string size. The approach used for the dictionary generation substantially impacts the compression ratio delivered by the algorithm. We can also compress data by changing the encoding of the output. Entropy coders, for example, try to reach the lower bound on the number of bits required to represent the symbols by leveraging the frequency of input patterns. Huffman coding and Lempel-Ziv are examples of entropy coders [11].

SMILES compression can be achieved by using state-of-the-art binary compression tools such as Bzip2¹, DEFLATE, and LZ77 [12], or by relying on tools for short string compression like SMAZ², SHOCO³, and FSST [13]. None of the previously mentioned tools meet our defined requirements: readable output, separable SMILES, and a shared dictionary. In the experimental results, we compare ZSMILES with Bzip2, as a representative of binary compressors file-based, and with FSST and SHOCO, as representative of small string compression.

Bzip2 compresses files using multiple layers ranging from the Burrows-Wheeler and move-to-front transform [14] to the Huffman coding [11]. Its compression ratio is high compared to other algorithms at the expense of time and computation [12]. Bzip2 compression is stateful, so to decompress one part of the file, it also has to decompress the previous part; thus, random access is impossible. Therefore, to meet our use case requirements, we can use Bzip2 to compress SMILES files line by line, which makes SMILES compression inefficient since the input string needs to be larger to reach a good compression ratio. Nonetheless, the compressed SMILES with Bzip2 are in binary format, so they are not human-readable.

SHOCO⁴ is a compression library for short strings. It provides a way to generate a custom dictionary based on the application domain. Since it is an entropy encoder, it provides

¹Website: <https://sourceware.org/bzip2/>

²Website: <https://github.com/antirez/smaz>

³Website: <https://ed-von-schleck.github.io/shoco/>

⁴Website: <https://ed-von-schleck.github.io/shoco/>

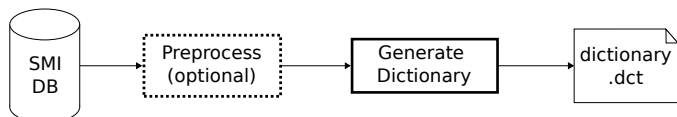


Fig. 2: Graphical representation of the SMILES pre-processing step.

a non-readable and non-random access output. FSST instead offers a good compression ratio on small strings within a dataset while allowing random access to the file. FSST uses a static symbol table defined from a small chunk of data from the input file. Since the table is static, it is immutable during the compression and decompression. FSST constructs a symbol table for each input; thus, the dictionary is input-dependent. In addition to that, FSST’s output is non-readable since the dictionary uses non-printable ASCII symbols, which can cause problems with third-party tools.

Another work available in the literature which leverages domain knowledge to achieve better SMILES compression through data preprocessing was presented by Gupta et al. [15]. It targets compression using file-wide binary compression tools such as Bzip2. This approach does not apply to our case because the compressed file cannot be accessed line-by-line, which does not guarantee the possibility of random access.

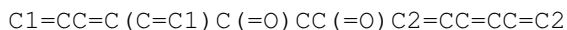
IV. METHODOLOGY

This section introduces the proposed ZSMILES approach. The main idea is to use a dictionary that associates a common string pattern to an ASCII character. To compress SMILES, we substitute all the dictionary patterns in SMILES with the related characters. We escape in output any characters in the input SMILES that cannot be encoded using the dictionary. The first two sections explain two optimizations that hinge on domain knowledge to improve the quality of the dictionary. Then, Section IV-C describes how we generate the dictionary in more detail, and Section IV-D reports the compression and decompression algorithm. Ultimately, we detail the CUDA accelerated SMILES compression approach in Section IV-E.

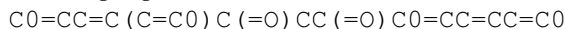
A. Preprocessing

SMILES pre-processing aims to increase the probability of finding common patterns. SMILES use numbers (IDs) to identify the opening and closing of rings. Since ring IDs are not reused, it became more difficult to identify common patterns. We propose pre-processing the inputs to increase the reuse of ring enumerations and the probability of finding common patterns. We can reuse ring IDs in 2 ways: innermost or outermost. If multiple rings are nested within each other, we can give the lower ID to the innermost or the outermost. We chose the innermost approach because the simplest and most common rings are the inner ones, which have the smaller ID values.

Take as an example the Dibenzoylmethane SMILES representation:



this SMILES have two rings, $C1=CC=C(C=C1)$ and $C2=CC=CC=C2$, which have a similar prefix, but with different IDs. Thus, compressing both of them would require the dictionary to have two entries representing the same prefix. Instead, if pre-processed, the SMILES becomes:



which now enables the compression of the SMILES by leveraging the occurrences of $C0=CC=C$. It is worth noting that the pre-processed SMILES remain valid.

B. Dictionary Pre-population

If the input SMILES has a pattern not included in the generated dictionary, we must escape all its characters, doubling the required size. Usually, known approaches mitigate this issue by using compression algorithms that build a dictionary tailored to the target input file. However, we would like to use a shared dictionary in our use case, increasing the chances of missing a pattern. For this reason, we can pre-populate the dictionary with the printable ASCII characters used by the SMILES format, thus avoiding escaping. For example, we use the character @ for chiral specification, the character / for stereoisomers specification, or the character # for triple bonds. This conservative choice reduces the number of characters representing patterns to the extended ASCII characters. However, when the input SMILES is compliant with the format, we have the guarantee that the compressed file does not require more storage than the input one. In the experimental results, we measured the compression ratio increment of this optimization.

C. Dictionary Generation

Figure 2 shows how ZSMILES generates the dictionary D , from an input training set of SMILES. The training set of SMILES is parsed to find a group of T recurrent substrings of the inputs that provide the highest coverage of the inputs. The coverage measures how much of the input is covered by the available substrings. This problem can, however, be seen as a knapsack problem [16], whose complexity is NP-complete.

Algorithm 1 is the pseudocode of how dictionaries are generated. The input is the training set of SMILES, while the output is the dictionary. The input is parsed by searching for unique substrings with length in the interval $[L_{min}, L_{max}]$. We used a maximum substring length of L_{max} because longer substrings require a lot of time to generate the dictionary. L_{min} was set to 2 to speed up the dictionary generation based on the previous detailed dictionary pre-population. Line 1 initializes $rank$, which contains the rank of all substrings found, while Line 2 initializes D , the compression dictionary. We define the $rank$ of a pattern p at step t as the product between its occurrences in the input, times its normalized length $occ_{pattern} \times norm_{p,t}$, as in Equation 1:

$$rank_{p,t} = occ_p \times (l_p - overlap_{p,t}) \quad (1)$$

Where $norm_{p,t}$ of pattern p at step t is defined as the difference between the length l_p of p and the overlap with patterns selected in the previous iteration of the loop in Line 8.

Algorithm 1: Dictionary generation algorithm.

Input: SMILES \leftarrow training set
Output: $D \leftarrow$ substring dictionary
Parameters: $L_{min} \leftarrow$ Minimum substring length
 $L_{max} \leftarrow$ Maximum substring length
 $C \leftarrow$ Initial dictionary values
 $T \leftarrow$ Dictionary size

```
1 rank  $\leftarrow$  {};  
2  $D \leftarrow C$ ;  
3 foreach line  $\in$  SMILES do  
4   foreach  $s \in$  line :  $L_{min} \leq |s| \leq L_{max}$  do  
5     rank[s] =  $\begin{cases} \text{rank}[s] + 1 & \text{if } s \in \text{rank}, \\ 1 & \text{otherwise} \end{cases}$   
6   end  
7 end  
8 foreach  $t \in \{0, 1, \dots, T\}$  do  
9   new_word =  $\text{argmax}_{\text{rank}(s)} s \in \text{rank}$ ;  
10   $D \leftarrow D \cup \{\text{new\_word}\}$ ;  
11  rank  $\leftarrow$  rank  $\setminus$  new_word;  
12  foreach  $s \in$  rank do  
13    update_rank(s);  
14  end  
15 end
```

From Line 3 to Line 7, count the occurrences of all substrings in the input. From Line 8 to Line 15, we populate D (size T), and at each iteration, we get the pattern that has the higher rank. Then, all other ranks are updated based on the pattern selected at that step.

D. ZSMILES Execution Flow

Figure 3 exemplify ZSMILES compression and decompression process. The flow in the upper part of the diagram illustrates the compression process: ZSMILES optionally preprocess the input before compressing and storing it. The lower part of the diagram instead reports the decompression process of ZSMILES, which works backward: from storage, compressed SMILES are decompressed and optionally post-processed. The dictionary is soft-coded in the ZSMILES executable, so we cannot change it once ZSMILES is compiled.

1) *Compression Algorithm:* We formulate the compression algorithm as an optimization problem. The input is a single line from a SMILES file. The output is the compressed SMILES. The problem is identifying the sequence of substrings in D that gives the best compression ratio of the input. Dictionary D is represented by a trie [17] to do pattern matching on the input symbols. We use the Dijkstra algorithm to get the shortest path from the first char in SMILES up to the last one. To apply Dijkstra, we need to construct a graph G of the SMILES: each node n is an input character, and an edge (n_1, n_2) represents a substring in D which begins with character n_1 and ends with character n_2 . The algorithm scans the graph G to match patterns from the trie of D . At each iteration i , the Dijkstra algorithm computes the best way

to compress the input from character n_i up to the input end, based on a cost function. Each match (represented by a graph edge) costs one if it is from a symbol in D . Otherwise, the cost is two if the symbol has to be escaped with a trailing *space* because there was no match in D . The last iteration calculates the shortest path to compress the input SMILES and traverses it from the beginning while printing the symbol of the corresponding pattern match on the output.

2) *Decompression:* The decompression algorithm is instead straightforward. During decompression, we use the dictionary D as a lookup table. For each symbol in each compressed SMILES, we perform a lookup in D and print out its expansion. If, instead, the value is a *space*, due to the escaping, we go to the next symbol in the input and print this symbol directly.

E. CUDA Implementation

To reduce ZSMILES compression (and decompression) overhead, we have developed a CUDA implementation, which allows us to exploit the computational power offered by NVIDIA GPUs. Given an input set of SMILES, the compression or decompression is split among groups of CUDA threads, called *blocks*, which build up a CUDA *grid*.

For compression, each block compresses a SMILES. Each block's thread looks at different input SMILES's characters: for each dictionary element, the thread checks if the correspondent substrings can be matched in the input, starting from that character. In this way, the block constructs a graph representation of the SMILES: a substring match creates an edge, which connects the substring's first character node with the node of the last character. We apply weights on the edges using a cost function in the same way as described in Section IV-D1. Once a block has a graph representation of the input SMILES, it scans the graph backward by applying Dijkstra. The shortest path identified by Dijkstra gives the best compression ratio of the input SMILES.

For decompression, each block decompresses a SMILES. Each block's thread performs a lookup into the dictionary using an input SMILES character. Therefore, each thread knows the decompressed string's dimension for each input character. Finally, block threads share how many characters they must write in output in order to know where to start the writing.

We have set blocks to have the same dimension as a CUDA warp. A warp is a set of 32 CUDA threads with specific scheduling properties. In our case, we rely on warps' synchronization and shuffle CUDA operations for performance reasons.

V. RESULTS

In this section, we discuss and report the experimental results of ZSMILES. We provide details of the setup and data used in Section V-A. Then, we analyze and discuss ZSMILES based on the compression ratios (Section V-B) and the performance (Section V-C).

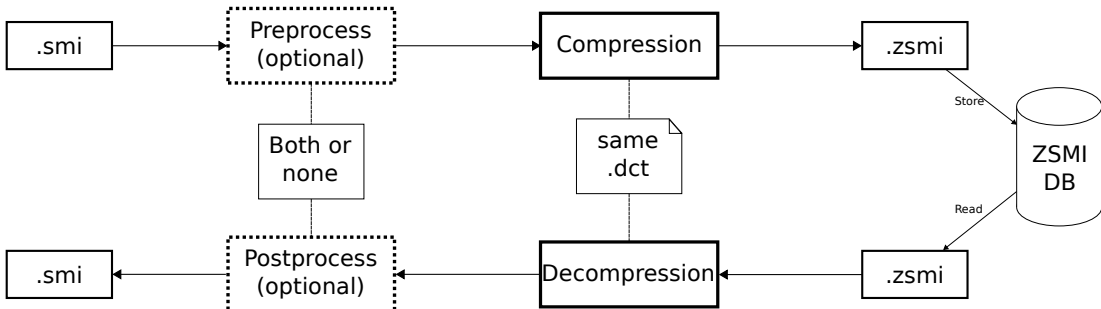


Fig. 3: Graphical representation of the compression and decompression process in ZSMILES.

A. Experimental Setup

Our experiment used a machine equipped with an AMD EPYC 7282 16-core processor, 64 GB of RAM, and two NVIDIA A100 graphics cards.

Since the dictionary generation has to be input-independent, we need a heterogeneous set of databases. Thus, we used three datasets that we consider representative: two of them, GDB-17 [18] and MEDiate database [19] are public, while one of them is a set of SMILES from a real case virtual screening execution [2], that we call EXSCALATE. GDB-17 contains 166 billion small organic molecules. MEDiate is a dataset of ligands from commercial compounds to natural products. We also used a MIXED dataset to construct the dictionary, combining the first one million ligands from each dataset. Dataset and additional material will be available on a public repository on GitHub⁵.

B. Compression Ratio

In this section, we report the results of different experiments on the ZSMILES compression ratio.

Dictionary Optimizations. The first experiment evaluates the ZSMILES compression ratio with different dictionary generation mechanisms. In particular, we evaluate the impact of what has been proposed in Section IV-A and in Section IV-B. To train the dictionaries used in Table I, we used a sample of random 50000 SMILES from the mixed dataset, the same one used to test the compression ratio in the experiment. The experiment collects ZSMILES compression ratios with all combinations of optimizations. We compare the compression ratio of ZSMILES when the dictionary is pre-populated with all the printable ASCII characters, with the characters of the SMILES alphabet, or with no characters at all. We expect a better compression ratio when the dictionary is initialized with a subset of common characters in SMILES.

Table I reports the experiment’s results. The first column indicates whether the pre-processing has been done on the input data. The second column reports the set of ASCII characters used to pre-populate the dictionary. From Table I, we can see how, in all cases, the reuse of ring ID has improved the compression ratio, and we can also see that pre-populating

the dictionary with the SMILES alphabet provides a better compression ratio, up to 0.29.

In the remaining part of this work, we apply pre-processing before generating all the dictionaries that are initialized with the SMILES alphabet. These are the proposed optimization techniques.

Cross-dictionary. The second experiment we conducted aims to find the best dataset to use for dictionary training. The experiment analyses the tradeoffs of using the same shared dictionary for any input set of SMILES. We have evaluated ZSMILES compression ratios with dictionaries trained on each available dataset against all others. We expect the compression to be worse when the training set consists of similar SMILES and better when we try to compress a dataset using a dictionary trained on the same one.

In Table II are the compression ratios of ZSMILES when the dictionary is trained on the dataset in the first column and tested on the dataset in the first row. When the dictionary is trained on GDB-17, MEDiate, or EXSCALATE, the average compression ratios obtained by compressing other datasets are 0.52, 0.34, and 0.39, respectively. The compression ratio can vary depending on the chosen training dataset. The dictionary generated on GDB-17 performs poorly on other datasets, which indicates that GDB-17 consists of homogenous SMILES. As expected, using a MIXED dataset gave an

TABLE I: Compression ratios of ZSMILES using different dictionaries.

Pre-processing	Pre-population	Compression Ratio
Yes	Printable	0.32
No	Printable	0.35
Yes	SMILES alphabet	0.29
No	SMILES alphabet	0.32
Yes	None	0.33
No	None	0.35

TABLE II: Compression ratios of ZSMILES using cross-dictionaries.

Train \ Test	Test			
	GDB-17	MEDIATE	EXSCALATE	MIXED
GDB-17	0.33	0.60	0.60	0.55
MEDIATE	0.46	0.29	0.29	0.35
EXSCALATE	0.52	0.36	0.31	0.38
MIXED	0.39	0.33	0.30	0.29

⁵Website: <https://github.com/elvispolimi/zsmiles>

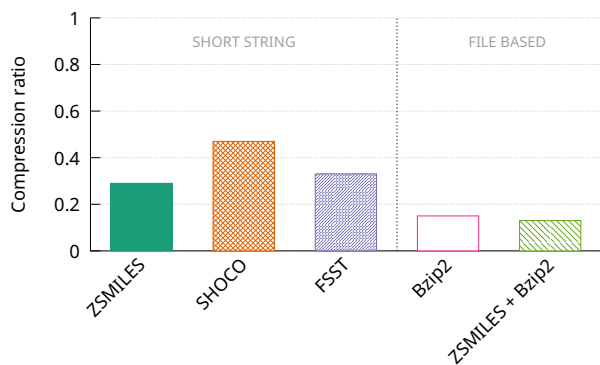


Fig. 4: Compression ratios of different tools on a mixed dataset, comparing both short-string and file-based methods.

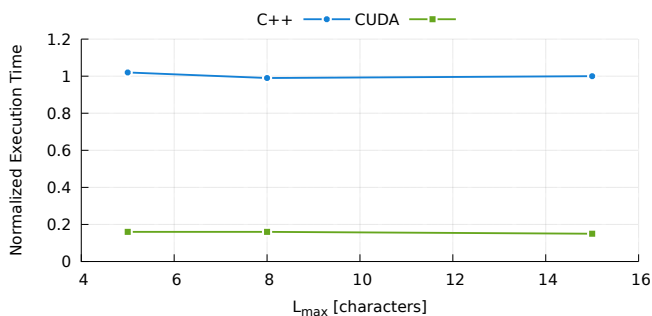
average compression ratio of 0.32, better than all the others. For this reason, we chose the MIXED dataset for the following experiments.

Tools Comparison. Finally, we compare the performance of ZSMILES with other state-of-the-art solutions: Bzip2 and FSST, to highlight the benefits of ZSMILES’s preprocessing optimizations. We compare these tools by measuring the compression ratio achieved. We use this as a test dataset, the MIXED one. We compressed the same MIXED dataset with each tool and trained ZSMILES’s dictionary on the same dataset: we made this decision because FSST constructs a static dictionary for each input, thus allowing us to compare the two approaches fairly. We expected the binary compression tool to yield the best compression ratio.

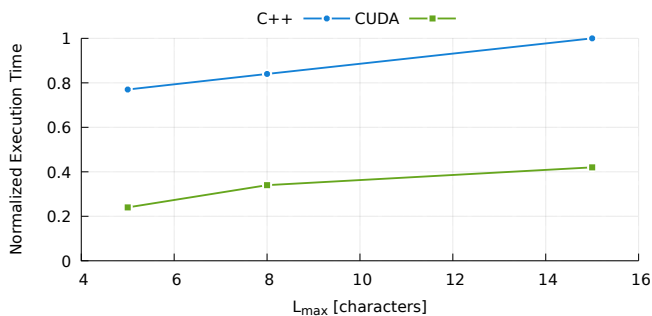
Figure 4 compares the compression ratio of ZSMILES with FSST, SHOCO, and Bzip2 on a MIXED dataset. The proposed domain-specific optimizations have been employed only by ZSMILES in this comparison. On the x-axis is the tool’s name used, while on the y-axis is the achieved compression ratio. Figure 4 indicates that ZSMILES can provide a good compression ratio while still producing a readable output. As expected, Bzip2 is the one that performs better, but it does not allow random access or reading of the output since it is binary and the compression is stateful. ZSMILES performs better than FSST when compressing on the same dataset used for dictionary training. Using Bzip2 on ZSMILES’s output can save even more space as it further compresses the data. It demonstrates the benefits of the preprocessing data step in ZSMILES for the BZIP2 compression algorithm.

C. Performance

In this section, we analyze ZSMILES compression and decompression performance: we report ZSMILES execution time by varying the maximum pattern length L_{max} in compression and decompression and with different implementations (C++ and CUDA). C++ refers to the serial implementation targeting CPUs, while CUDA refers to the parallel one targeting NVIDIA’s GPUs. All execution times reported here consider the execution time of the entire ZSMILES application, and



(a) Compression performance.



(b) Decompression performance.

Fig. 5: ZSMILES normalized execution times of the C++ and CUDA implementation with different L_{max} values.

they are normalized to the execution time of the C++ one with the maximum value of L_{max} . We expect the CUDA version to be faster than the C++ one. We evaluate the execution time of ZSMILES C++ and CUDA version on the Mixed dataset, with different values of L_{max} : 5, 8, 15.

Figure 5 reports ZSMILES execution times. L_{max} values are on the x-axis, while ZSMILES normalized execution times are on the y-axis. Figure 5a shows ZSMILES’s normalized execution times in compression, while Figure 5a in decompression. As expected, the CUDA implementation is faster in compression and decompression than the C++ one. In particular, the CUDA version is only 7× faster in compression and 2× in decompression. We have investigated these speedups further and discovered that for these simple kernels, the bottlenecks are the read-and-write operations on storage: ZSMILES is memory-bound. Thus, additional C++ or CUDA optimizations have a reduced impact on performance.

VI. CONCLUSIONS

Virtual screening campaign requires the storage of large datasets of ligands. They involve trillions of ligands, which have a big storage footprint. Some examples in the literature [8] report the burden of storing large chemical spaces of molecules for extreme-scale campaigns. In this paper, we have proposed a methodology to store SMILES efficiently, called ZSMILES, by creating a shared dictionary through heuristics, data pre-processing, and dictionary pre-population. ZSMILES uses an efficient dictionary-based compression algorithm,

leveraging SMILES domain knowledge for optimizations and design. Given the use case, ZSMILES allows random access to the data and provides a readable output (ASCII format).

ZSMILES can reduce the space required for large datasets by providing a compression ratio of up to 0.29, showing an improvement over state of the art approach of $\times 1.13$.

We propose a C++ and CUDA implementation of ZSMILES, where the CUDA parallel one achieves a speedup of $7\times$ in compression and $2\times$ in decompression compared to the C++ serial implementation.

From experimental results, we have found ZSMILES compression and decompression overhead to be negligible since kernels are inherently memory-bound.

ACKNOWLEDGMENT

This project has received funding from EuroHPC-JU - the European High-Performance Computing Joint Undertaking - under grant agreement No 956137 (LIGATE). The JU receives support from the European Union's Horizon 2020 research and innovation program and Italy, Sweden, Austria, Czech Republic, and Switzerland.

REFERENCES

- [1] M. Schlander, K. Hernandez-Villafuerte, C.-Y. Cheng, J. Mestre-Ferrandiz, and M. Baumann, "How Much Does It Cost to Research and Develop a New Drug? A Systematic Review and Assessment," *PharmacoEconomics*, vol. 39, no. 11, pp. 1243–1269, Nov. 2021.
- [2] D. Gadioli, E. Vitali, F. Ficarelli, C. Latini, C. Manelfi, C. Talarico, C. Silvano, C. Cavazzoni, G. Palermo, and A. R. Beccari, "Exscalate: An extreme-scale virtual screening platform for drug discovery targeting polypharmacology to fight sars-cov-2," *IEEE Transactions on Emerging Topics in Computing*, vol. 11, no. 1, pp. 170–181, 2023.
- [3] G. K. Kiriiri, P. M. Njogu, and A. N. Mwangi, "Exploring different approaches to improve the success of drug discovery and development projects: A review," *Future Journal of Pharmaceutical Sciences*, vol. 6, no. 1, p. 27, Dec. 2020.
- [4] N. A. Murugan, A. Podobas, D. Gadioli, E. Vitali, G. Palermo, and S. Markidis, "A review on parallel virtual screening softwares for high-performance computers," *Pharmaceuticals*, vol. 15, no. 1, 2022. [Online]. Available: <https://www.mdpi.com/1424-8247/15/1/63>
- [5] B. Zhang, H. Li, K. Yu, and Z. Jin, "Molecular docking-based computational platform for high-throughput virtual screening," *CCF Transactions on High Performance Computing*, pp. 1–12, 2022.
- [6] E. H. B. Maia, L. C. Assis, T. A. de Oliveira, A. M. da Silva, and A. G. Taranto, "Structure-based virtual screening: From classical to artificial intelligence," *Frontiers in Chemistry*, vol. 8, Apr. 2020. [Online]. Available: <https://doi.org/10.3389/fchem.2020.00343>
- [7] S. LeGrand, A. Scheinberg, A. F. Tillack, M. Thavappiragasam, J. V. Vermaas, R. Agarwal, J. Larkin, D. Poole, D. Santos-Martins, L. Solis-Vasquez *et al.*, "Gpu-accelerated drug discovery with docking on the summit supercomputer: Porting, optimization, and application to covid-19 research," in *Proceedings of the 11th ACM international conference on bioinformatics, computational biology and health informatics*, 2020, pp. 1–10.
- [8] D. Gadioli, E. Vitali, F. Ficarelli, C. Latini, C. Manelfi, C. Talarico, C. Silvano, C. Cavazzoni, G. Palermo, and A. R. Beccari, "Exscalate: an extreme-scale virtual screening platform for drug discovery targeting polypharmacology to fight sars-cov-2," *IEEE Transactions on Emerging Topics in Computing*, vol. 11, no. 1, pp. 170–181, 2022.
- [9] D. Weininger, "Smiles, a chemical language and information system. 1. introduction to methodology and encoding rules," *Journal of Chemical Information and Computer Sciences*, vol. 28, no. 1, pp. 31–36, 1988. [Online]. Available: <https://doi.org/10.1021/ci00057a005>
- [10] T. Gagie and G. Manzini, *Dictionary-Based Data Compression*. Boston, MA: Springer US, 2008, pp. 236–240. [Online]. Available: https://doi.org/10.1007/978-0-387-30162-4_108
- [11] D. A. Huffman, "A method for the construction of minimum-redundancy codes," *Proceedings of the IRE*, vol. 40, no. 9, pp. 1098–1101, 1952.
- [12] A. Gupta, A. Bansal, and V. Khanduja, "Modern lossless compression techniques: Review, comparison and analysis," in *2017 Second International Conference on Electrical, Computer and Communication Technologies (ICECCT)*, 2017, pp. 1–8.
- [13] P. Boncz, T. Neumann, and V. Leis, "Fsst: fast random access string compression," *Proc. VLDB Endow.*, vol. 13, no. 12, p. 2649–2661, jul 2020. [Online]. Available: <https://doi.org/10.14778/3407790.3407851>
- [14] G. Manzini, "An analysis of the burrows—wheeler transform," *J. ACM*, vol. 48, no. 3, p. 407–430, may 2001. [Online]. Available: <https://doi.org/10.1145/382780.382782>
- [15] S. Scanlon and M. Ridley, *A Fully Reversible Data Transform Technique Enhancing Data Compression of SMILES Data*. Springer Berlin Heidelberg, 2013, p. 54–68. [Online]. Available: http://dx.doi.org/10.1007/978-3-642-40511-2_5
- [16] H. Kellerer, U. Pferschy, and D. Pisinger, *Introduction to NP-Completeness of Knapsack Problems*. Springer Berlin Heidelberg, 2004, p. 483–493. [Online]. Available: http://dx.doi.org/10.1007/978-3-540-24777-7_16
- [17] E. Fredkin, "Trie memory," *Commun. ACM*, vol. 3, no. 9, p. 490–499, sep 1960. [Online]. Available: <https://doi.org/10.1145/367390.367400>
- [18] L. Ruddigkeit, R. van Deursen, L. C. Blum, and J.-L. Reymond, "Enumeration of 166 billion organic small molecules in the chemical universe database gdb-17," *Journal of Chemical Information and Modeling*, vol. 52, no. 11, pp. 2864–2875, 2012, pMID: 23088335. [Online]. Available: <https://doi.org/10.1021/ci300415d>
- [19] G. Vistoli, C. Manelfi, C. Talarico, A. Fava, A. Warshel, I. V. Tetko, R. Apostolov, Y. Ye, C. Latini, F. Ficarelli, G. Palermo, D. Gadioli, E. Vitali, G. Varriale, V. Pisapia, M. Scaturro, S. Coletti, D. Gregori, D. Gruffat, E. Leija, S. Hessenauer, A. Delbianco, M. Allegretti, and A. R. Beccari, "Mediate - molecular docking at home: Turning collaborative simulations into therapeutic solutions," *Expert Opinion on Drug Discovery*, vol. 18, no. 8, p. 821–833, Jul. 2023. [Online]. Available: <http://dx.doi.org/10.1080/17460441.2023.2221025>