# LZ1 Sequential Implementation
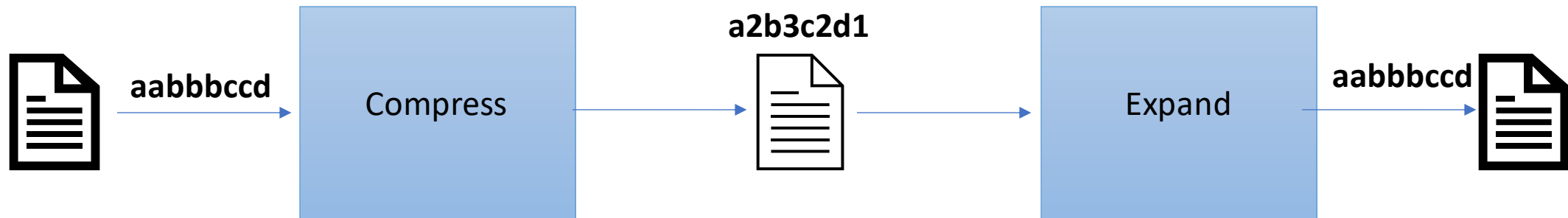
Presenters:
Rahul Mittal
Neel Prajapati

- In the early days of video gaming, fitting large games onto a single storage media was a challenge, making games expensive and limiting accessibility. But Capcom, the creators of Street Fighter 2, found an innovative solution - data compression.

- Capcom analysed and compressed the game's graphics, sound effects, and other data using an algorithm with a unique ability to recognize repeating patterns in the data, resulting in a visually stunning game that could fit on a single storage media, making it accessible to a wider audience.

- The algorithm used was LZ1, named after its inventors Lempel and Ziv. In this presentation, we will explore the mechanics of LZ1 and its use in compressing text files

# Data Compression?

- Data compression is the process of reducing the size of a computer file. Through an algorithm or a set of rules for carrying out an operation, computers can determine ways to shorten long strings of data and later reassemble them in a recognizable form upon retrieval.

# Basic Model for Data Compression

**aabbbccd** → Compress → **a2b3c2d1** → Expand → **aabbbccd**

# Data Has Repetitions!

- Data compression algorithms save space by exploiting the fact that most data files have many repeating patterns in them.

- For example, text data has certain character sequences that appear much more often than others

# Repetitions Are Beneficial

- Identifying repeating patterns makes it possible to represent them using less space.

- One approch to identification and space saving is to point to past occurrence of the pattern.

- LZ1 algorithm, that is the topic of this presentation, is one of the members of a well know family of algorithms that uses adaptive dictionaries as the reference mechanism to point to past occurrences.

# LZ1 Algorithm

```
1.  while (input is not empty)
2.      {
3.          get a reference (position, length) to longest match;
4.          if (length > 0)
5.          {
6.              output (position, length, next symbol);
7.              shift the window length+1 positions along;
8.          }else {
9.              output (0, 0, first symbol in the lookahead buffer);
10.             shift the window 1 character along;
11.         }
12.     }
```

# Compression Example

| | | | |
|---|---|---|---|
| | sir␣sid␣eastman␣ | $\Rightarrow$ | $(0,0,\text{"s"})$ |
| s | ir␣sid␣eastman␣e | $\Rightarrow$ | $(0,0,\text{"i"})$ |
| si | r␣sid␣eastman␣ea | $\Rightarrow$ | $(0,0,\text{"r"})$ |
| sir | ␣sid␣eastman␣eas | $\Rightarrow$ | $(0,0,\text{"␣"})$ |
| sir␣ | sid␣eastman␣easi | $\Rightarrow$ | $(4,2,\text{"d"})$ |
| sir␣sid | ␣eastman␣easily␣ | $\Rightarrow$ | $(4,1,\text{"e"})$ |
| sir␣sid␣e | astman␣easily␣te | $\Rightarrow$ | $(0,0,\text{"a"})$ |

Source: Data Compression-The Complete Reference

# An Advantage of Sliding Window

As the search window contains characters that are directly preceding the current characters being compressed, patterns repeating in only certain parts of the input are captured efficiently and effectively.

*wabbab̶wabbab̶wabbab̶wabbab̶woob̶woob̶woo*

(Source: Introduction to Data Compression)

| Encoder Output | Dictionary | |
|---|---|---|
| | Index | Entry |
| $\langle 0, C(w) \rangle$ | 1 | *w* |
| $\langle 0, C(a) \rangle$ | 2 | *a* |
| $\langle 0, C(b) \rangle$ | 3 | *b* |
| $\langle 3, C(a) \rangle$ | 4 | *ba* |
| $\langle 0, C(b̶) \rangle$ | 5 | *b̶* |
| $\langle 1, C(a) \rangle$ | 6 | *wa* |
| $\langle 3, C(b) \rangle$ | 7 | *bb* |
| $\langle 2, C(b̶) \rangle$ | 8 | *ab̶* |
| $\langle 6, C(b) \rangle$ | 9 | *wab* |
| $\langle 4, C(b̶) \rangle$ | 10 | *bab̶* |
| $\langle 9, C(b) \rangle$ | 11 | *wabb* |
| $\langle 8, C(w) \rangle$ | 12 | *ab̶w* |
| $\langle 0, C(o) \rangle$ | 13 | *o* |
| $\langle 13, C(b̶) \rangle$ | 14 | *ob̶* |
| $\langle 1, C(o) \rangle$ | 15 | *wo* |
| $\langle 14, C(w) \rangle$ | 16 | *ob̶w* |
| $\langle 13, C(o) \rangle$ | 17 | *oo* |

# UTF-8 Representation

**SPACE SAVER BECAME A NIGHTMARE**

## Code point ↔ UTF-8 conversion

| First code point | Last code point | Byte 1 | Byte 2 | Byte 3 | Byte 4 | Code points |
|---|---|---|---|---|---|---|
| U+0000 | U+007F | 0xxxxxxx | | | | 128 |
| U+0080 | U+07FF | 110xxxxx | 10xxxxxx | | | 1920 |
| U+0800 | U+FFFF | 1110xxxx | 10xxxxxx | 10xxxxxx | | [a]61440 |
| U+10000 | [b]U+10FFFF | 11110xxx | 10xxxxxx | 10xxxxxx | 10xxxxxx | 1048576 |

Source: Wikipedia

# Challenges Faced

String Indexing in C++

| 0 | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| अ | a | μ | # | ż | ʊ |

String Indexing in Julia

| 1110 0000 | 1010 0100 | 1000 0101 | 0011 1101 | 1100 0010 | 1011 0101 | 0001 0111 | 1110 0000 | 1010 0010 | 1010 0010 | 1100 0110 | 1011 0001 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | **2** | **3** | 4 | 5 | **6** | 7 | 8 | **9** | **10** | 11 | **12** |

```
1    a="ГДЄ(Е)Є"
2    a[2]
```

PROBLEMS    OUTPUT    DEBUG CONSOLE    TERMINAL

ERROR: LoadError: StringIndexError: invalid index [2], valid nearby indices [1]=>'Г', [3]=>'Д'

# Julia String Indices are mapped to bytes.

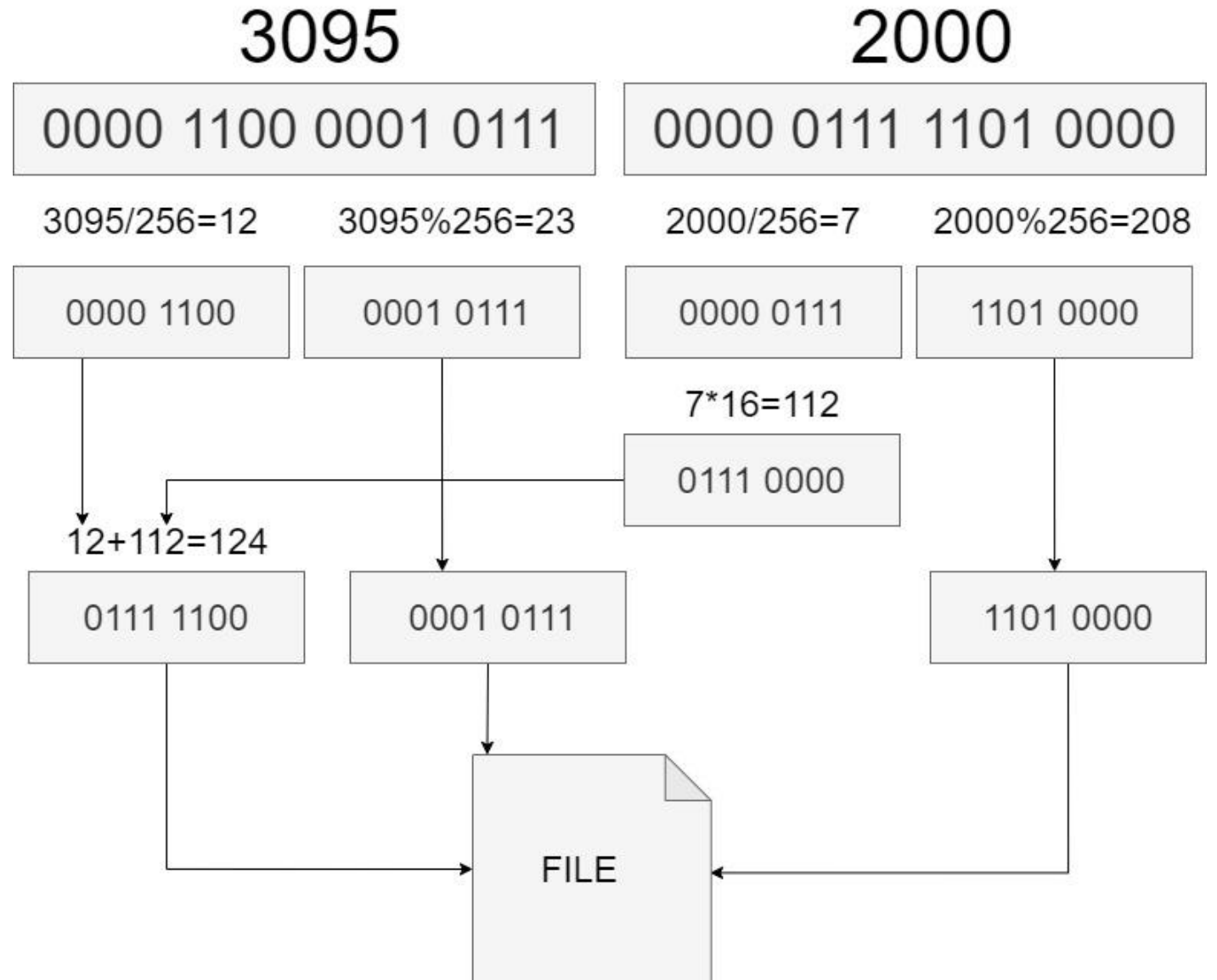# CODE COMPARISON

Restricted to ASCII in UTF-8

```
for i in target-current_search_buffer_length:target-1
```

For Unicode

```
eachindex_correction=prevind(input,target_index,current_search_buffer_length)-1
for index in eachindex(input[prevind(input,target_index,current_search_buffer_length):prevind(input,target_index)])
```
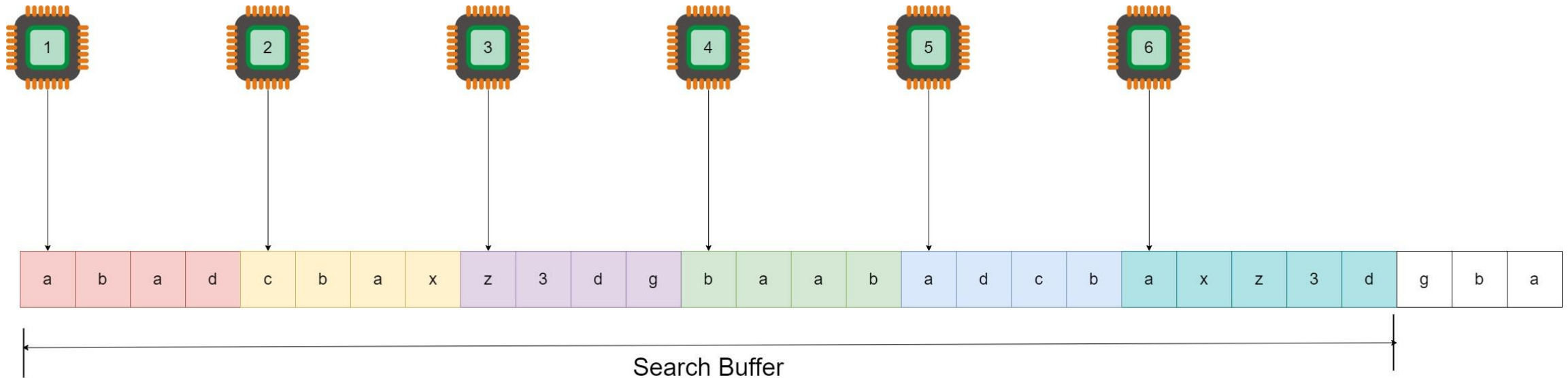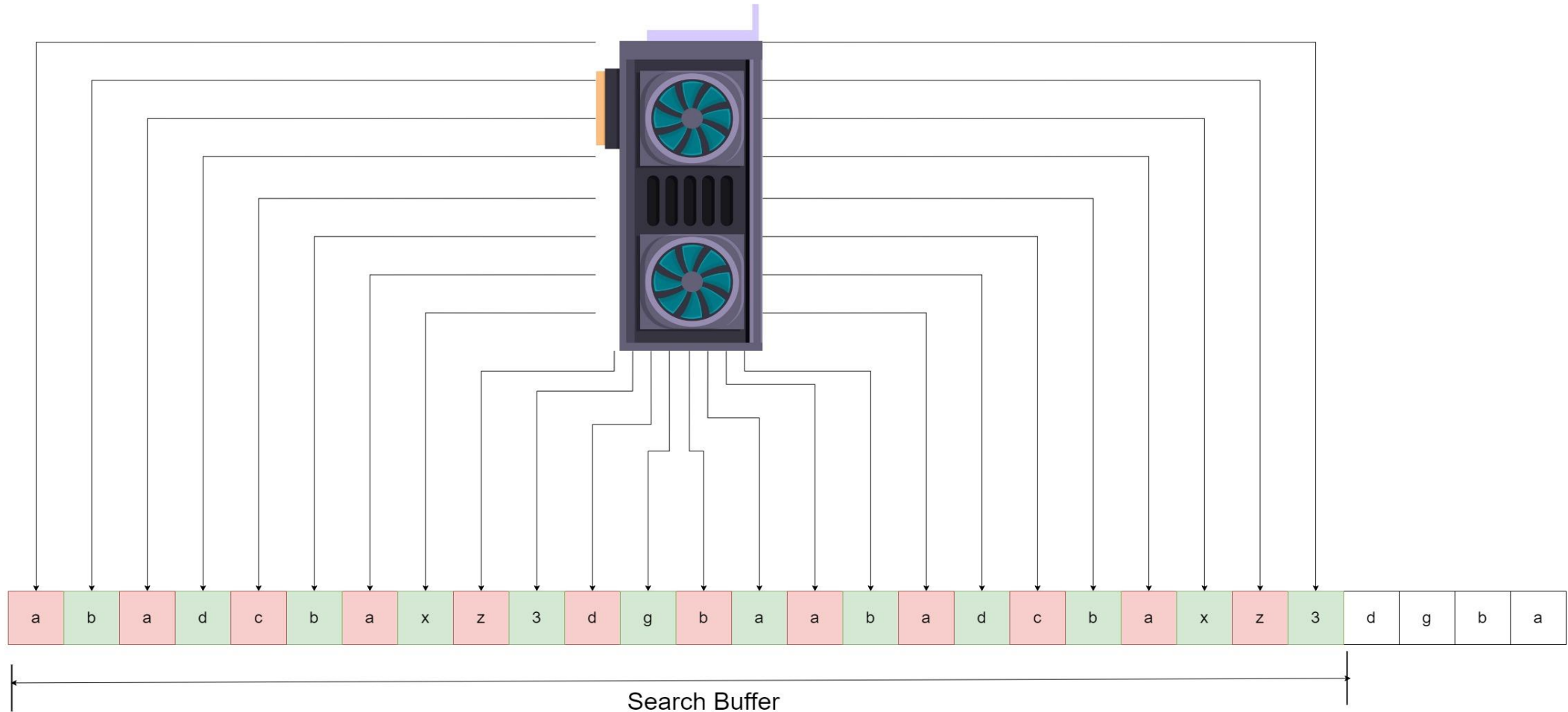
# Potential Opportunities for Parallel Implementations

## 1.Parallelizing the search phase on CPU:

## 2. Parallelizing the search phase on GPU:



Search Buffer

# 3. Dividing the File and Compressing in Parallel

| 1 | a | b | a | d | c | b | a | x | z | 3 | d | g | b | a | a | b | a | d | c | b | a | x | z | 3 | d | g | b | a |

Search Buffer

| 2 | a | b | a | d | c | b | a | x | z | 3 | d | g | b | a | a | b | a | d | c | b | a | x | z | 3 | d | g | b | a |

Search Buffer

| 3 | a | b | a | d | c | b | a | x | z | 3 | d | g | b | a | a | b | a | d | c | b | a | x | z | 3 | d | g | b | a |

Search Buffer

| 4 | a | b | a | d | c | b | a | x | z | 3 | d | g | b | a | a | b | a | d | c | b | a | x | z | 3 | d | g | b | a |

Search Buffer