

Suffix array construction

The “skew” algorithm

Sorting suffixes

mississippi\$
ississippi\$
ssissippi\$
sissippi\$
issippi\$
ssippi\$
sippi\$
ippi\$
ppi\$
pi\$
i\$
\$

Sorting suffixes...



\$
i\$
ippi\$
issippi\$
ississippi\$
mississippi\$
pi\$
ppi\$
sippi\$
sissippi\$
ssippi\$
ssissippi\$

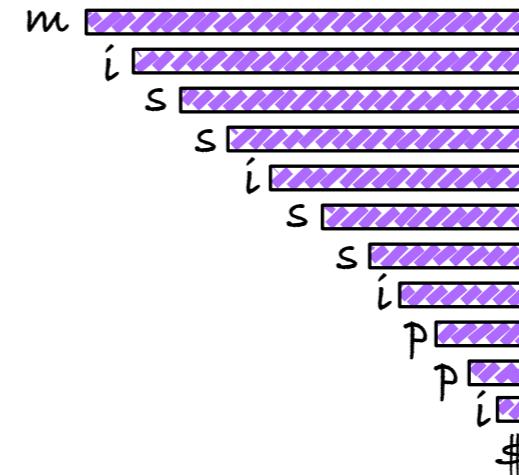
How long would it take with a standard sort algorithm (like quick sort?)

$O(n \log n)$ comparisons; $O(n)$ per comparison. $O(n^2 \log n)$ in total.

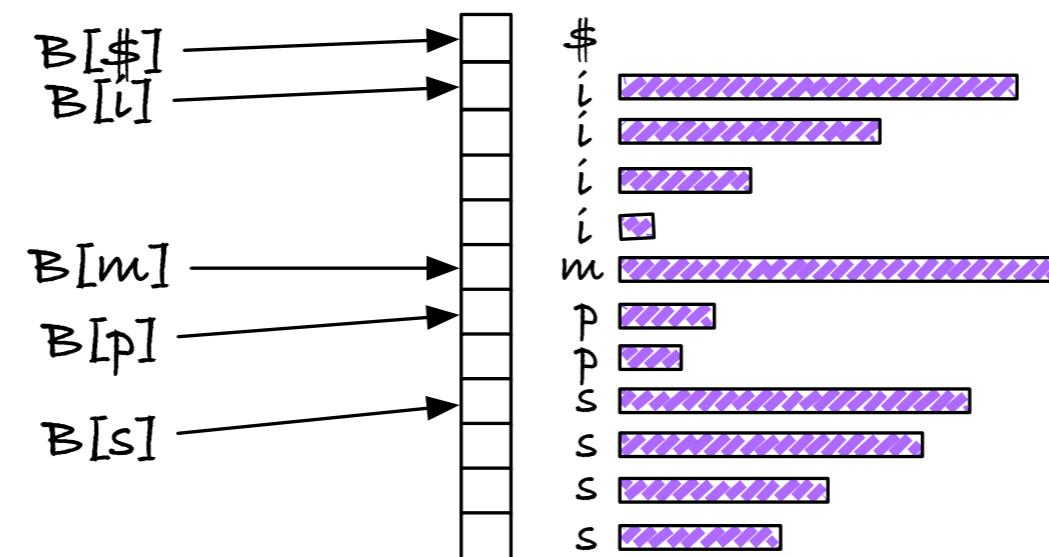
Bucket sort

Counts: \$ i m p s
 1 4 1 2 4

Buckets:
(acc sum of counts) \$ i m p s
 0 1 5 6 8

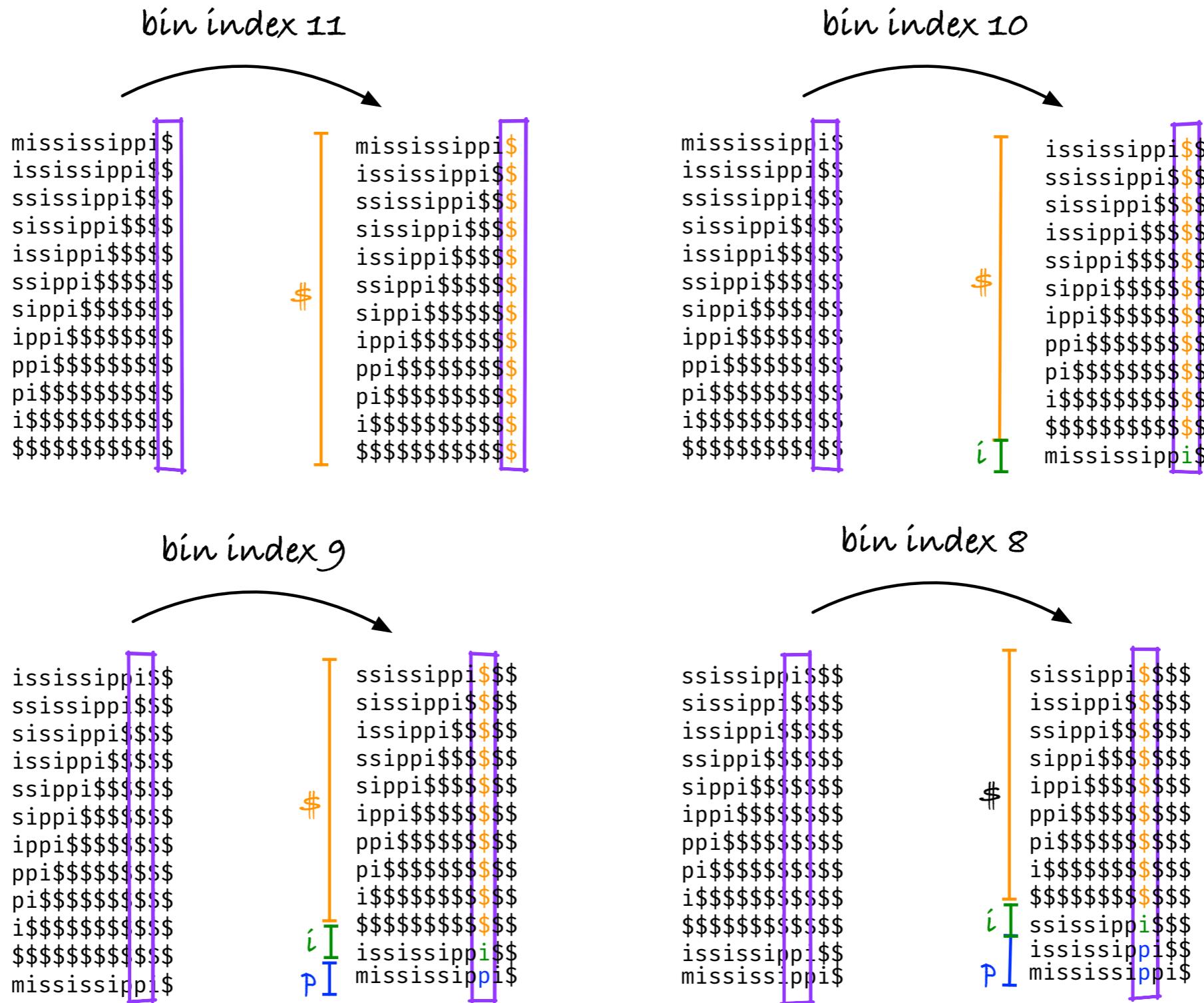


$$A[B[x[i]]++]=i$$

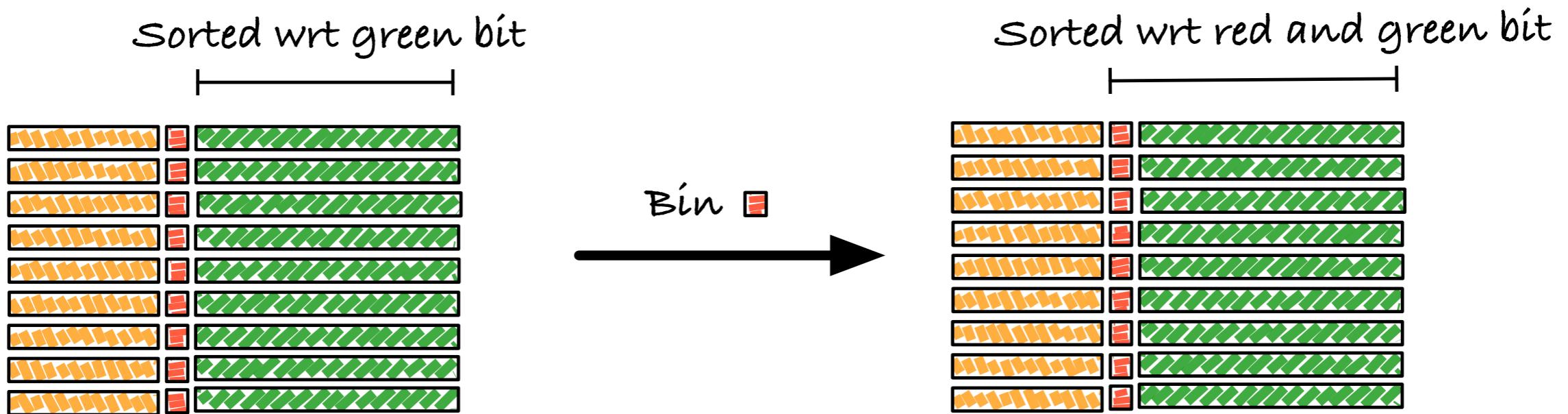


Running time: $O(n + \sigma) = O(n)$

Radix sort



Radix sort

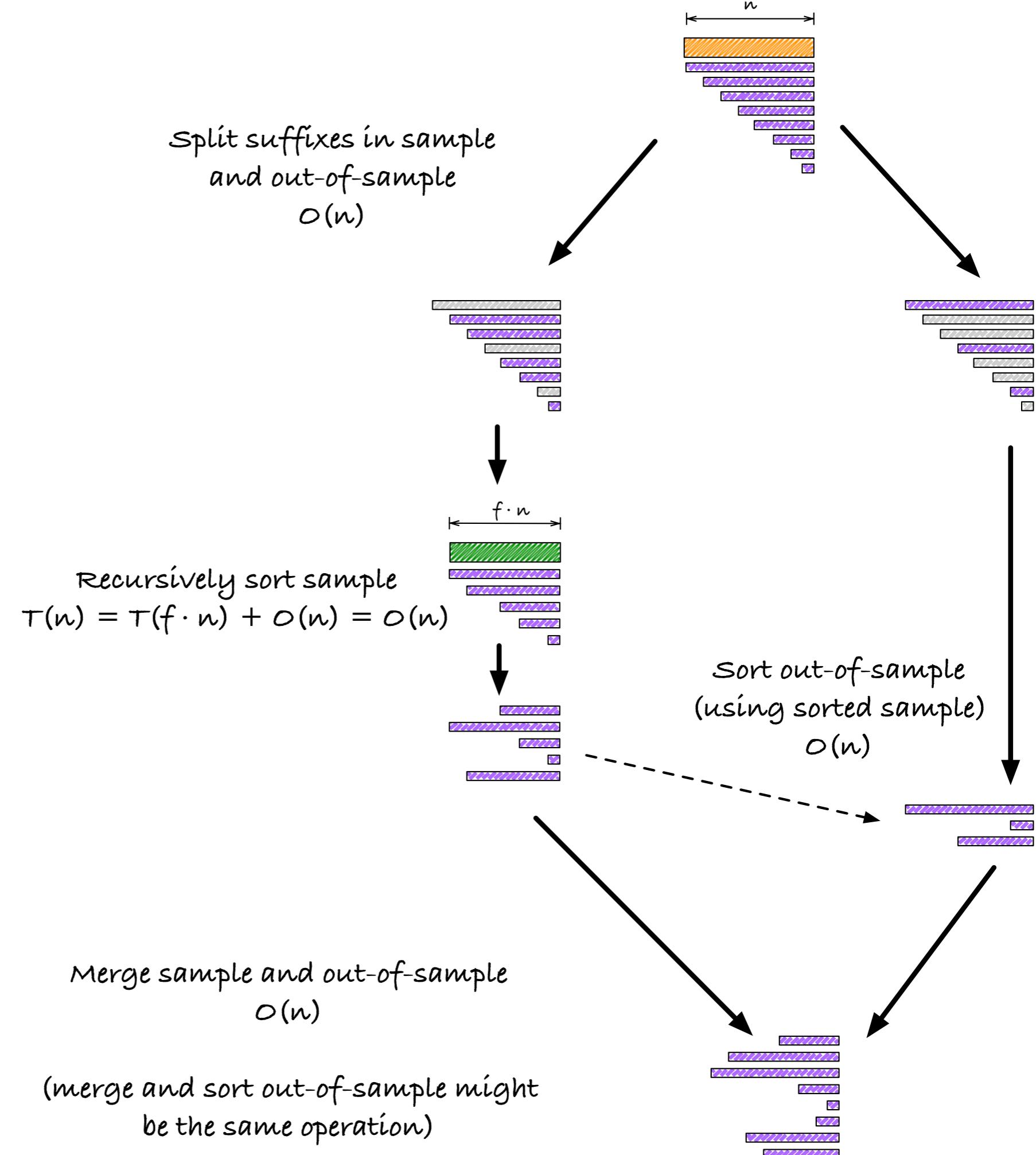


$O(n)$ per bucket sort.
 $O(n^2)$ for full radix sort.

Radix sort

- You can also radix sort from left to right by handling bins recursively
- Then you can terminate early when bins have size one
- but there is a little more overhead to it
- We won't be radix sorting long strings, so you can just use the simplest approach when you need to sort...

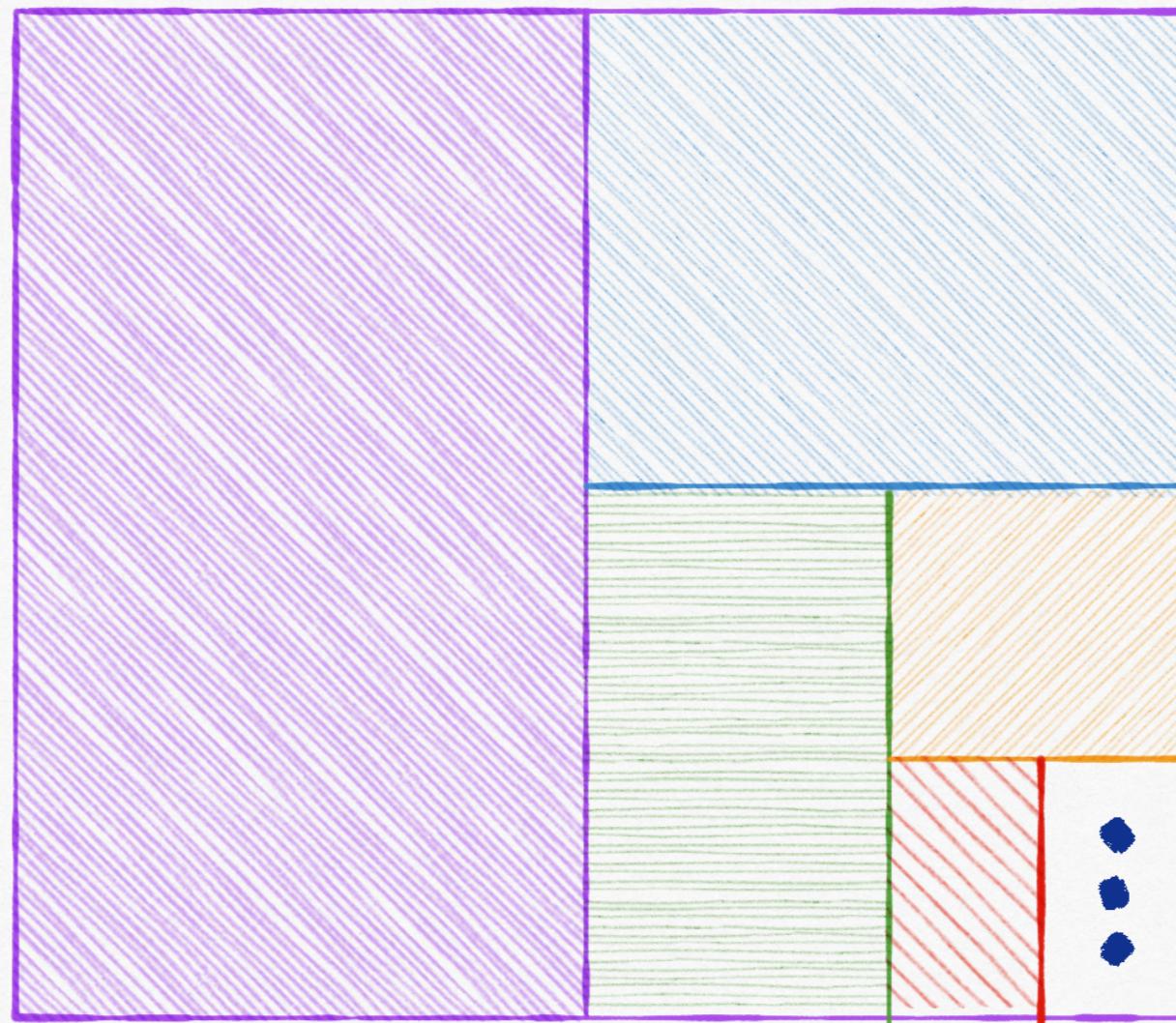
Sampling approach (to get to $O(n)$)



Recurrence

$$T(n) = T(f \cdot n) + O(n)$$

$$T(n) \in O\left(n \sum_{i=0}^{\infty} f^i\right)$$



$$n + n/2 + n/4 + n/8 + n/16 + \dots \leq 2n$$

Recurrence

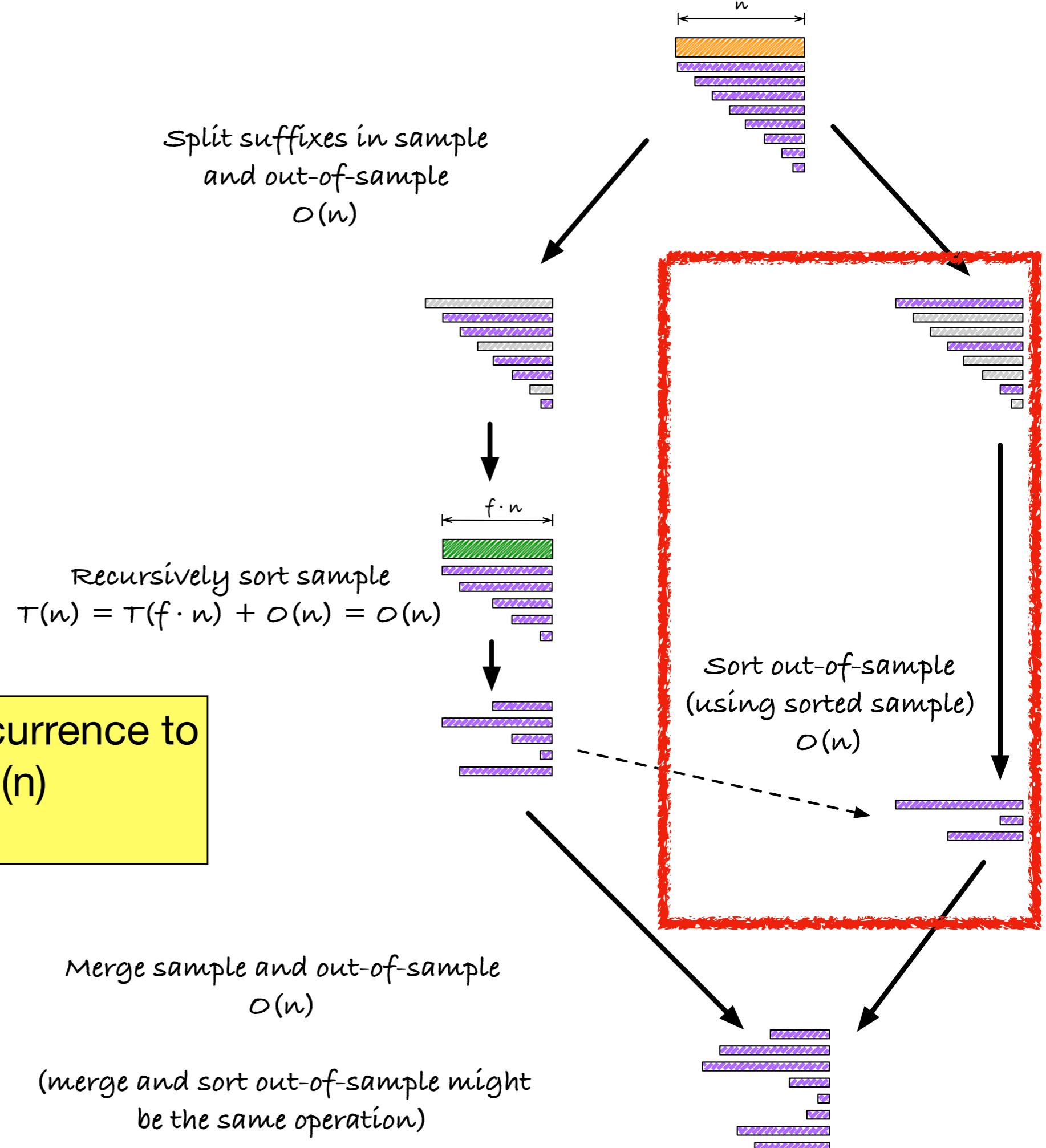
$$T(n) = T(f \cdot n) + O(n)$$

$$T(n) \in O\left(n \sum_{i=0}^{\infty} f^i\right)$$

$$\sum_{k=0}^{\infty} f^k = \frac{1}{1-f} \text{ for } 0 \leq f < 1$$

Why don't we sort out-of-sample recursively?

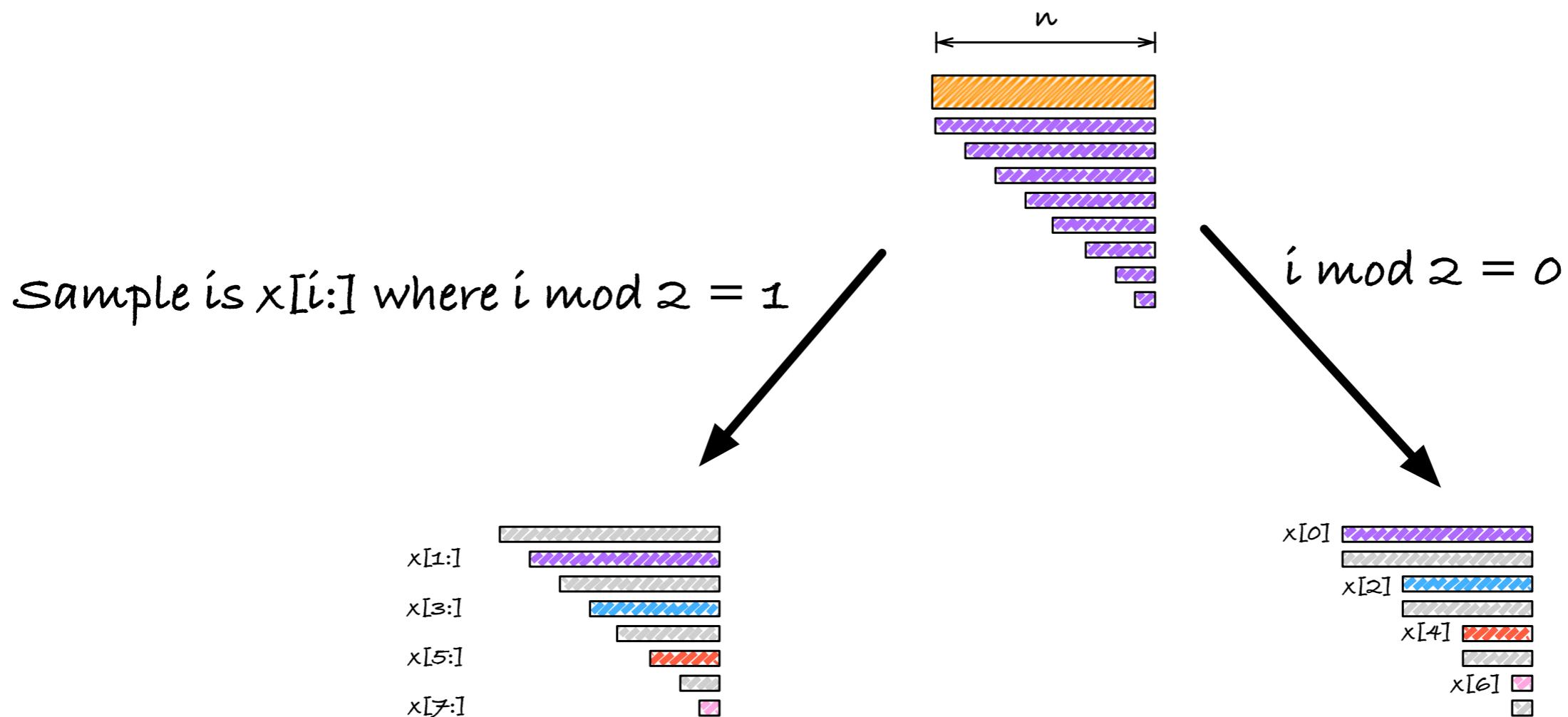
That would change the recurrence to
 $T(n) = T(f \cdot n) + T((1-f) \cdot n) + O(n)$
which isn't in $O(n)$.



First attempt

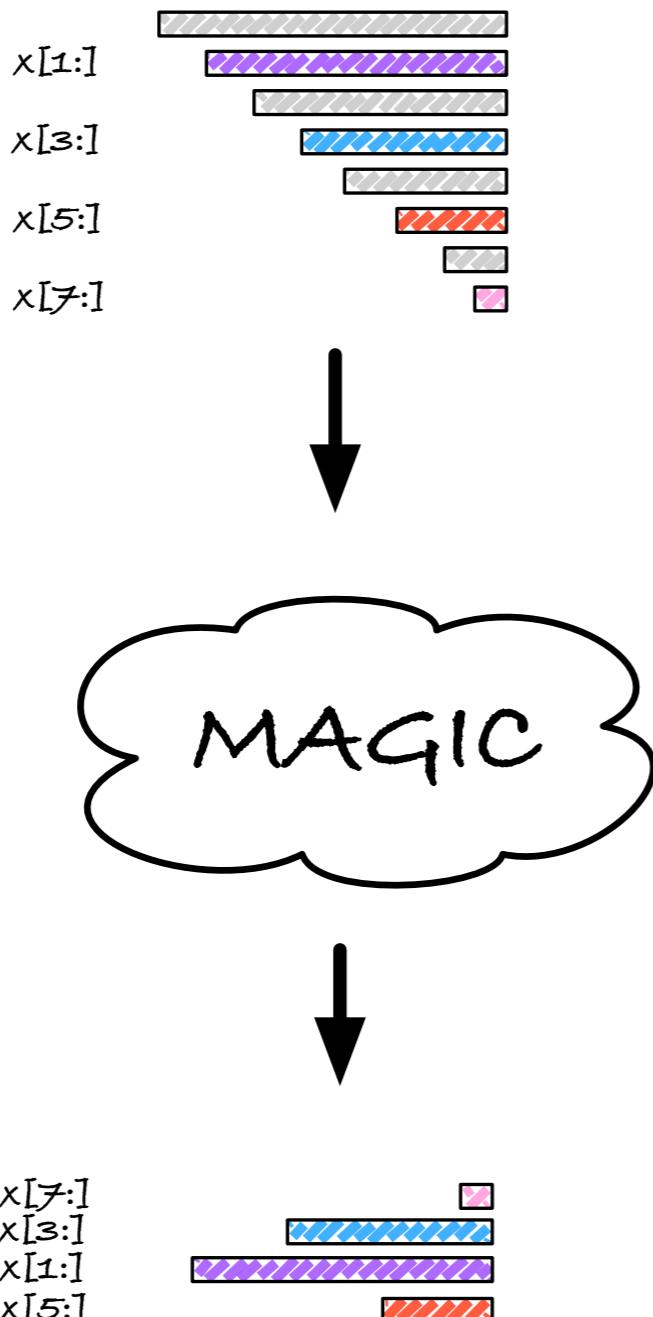
- The first attempt is the obvious choice for selecting samples
- It won't work, but the way it doesn't work will tell us why the skew algorithm doesn't do the obvious (and why it does what it does)

First attempt



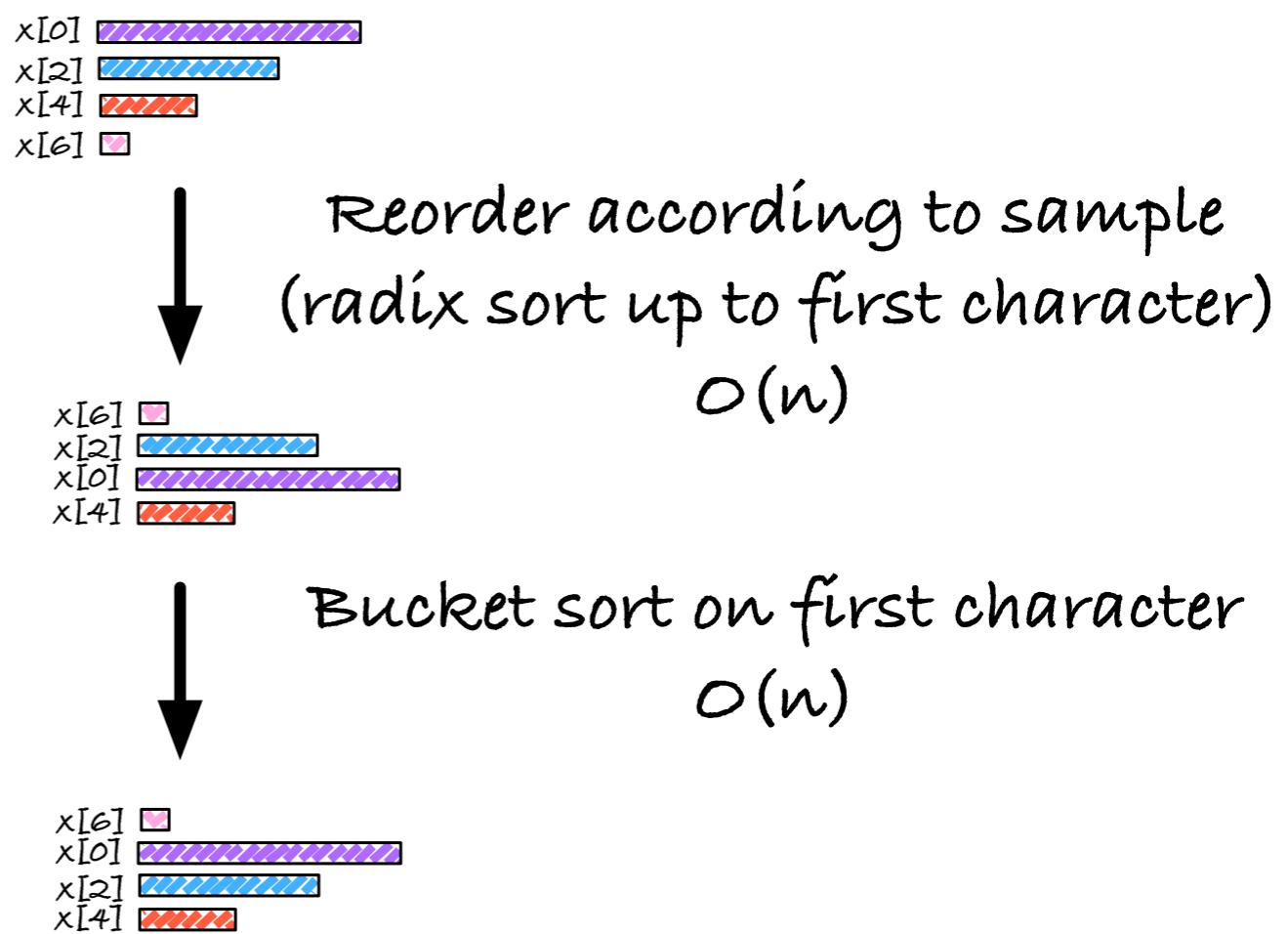
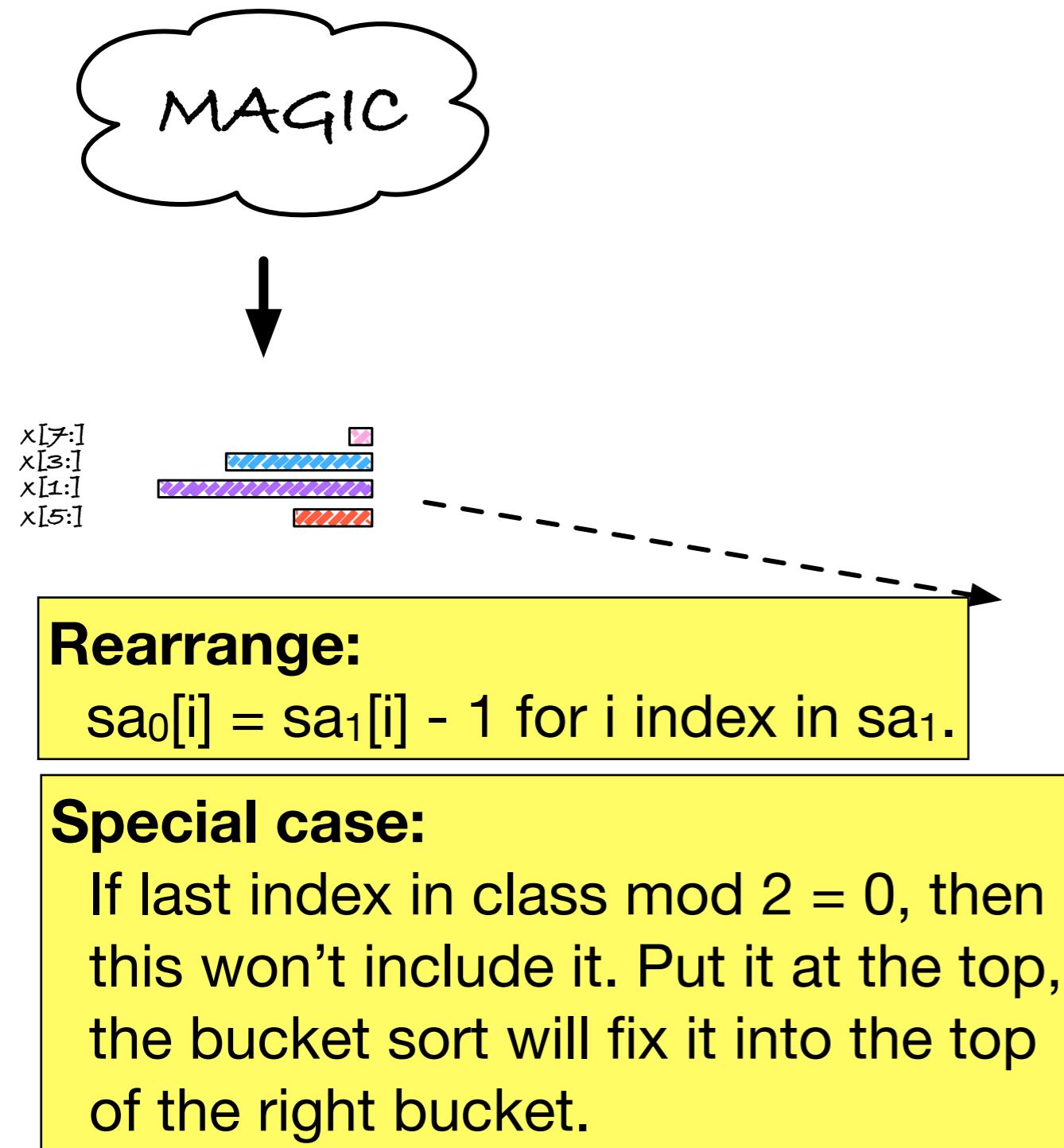
Easily done in $O(n)$ since it is just a scan
through the input string...

First attempt



You can get this to work, but I won't show it, because we will run into a problem shortly. I will show you how to sort the sample in the fixed version instead.

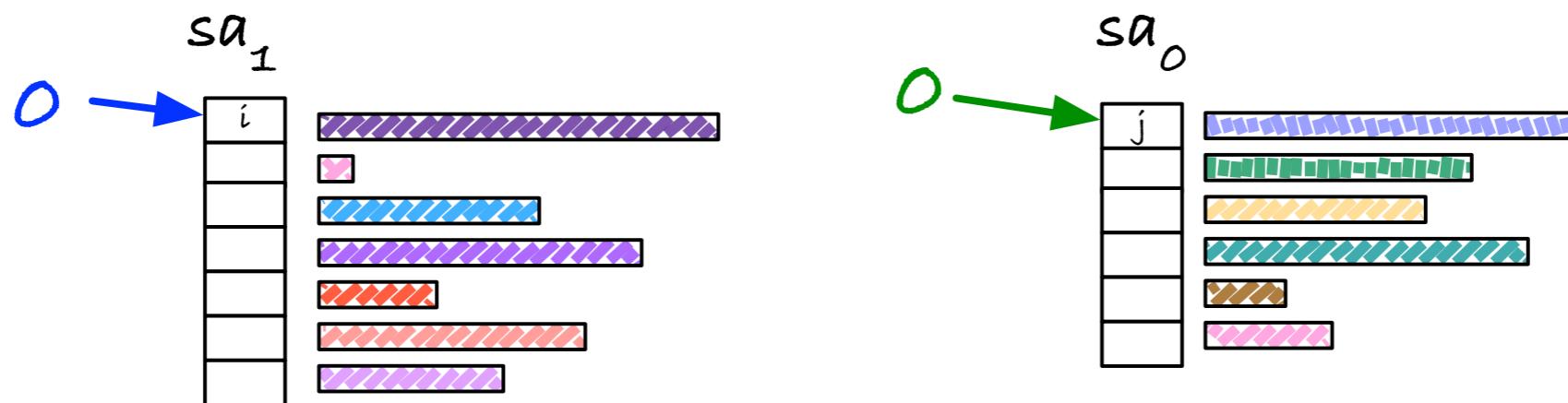
First attempt



First attempt

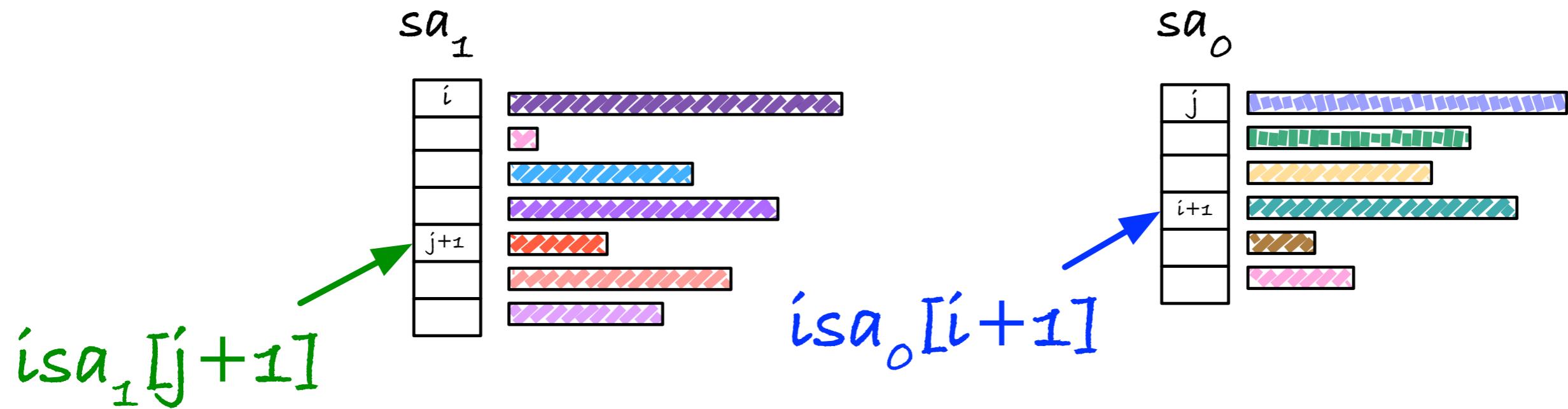
We can merge two arrays/lists in $O(c \cdot n)$ where c is the comparison time.

Can we compare these strings in $O(1)$?



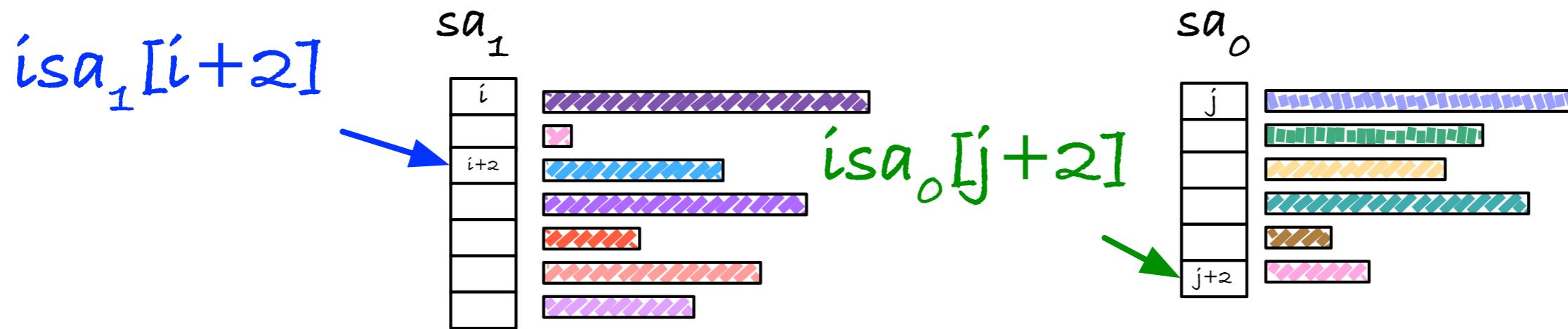
```
def less(x[i:], x[j]):  
    if x[i] < x[j]: return True  
    if x[i] > x[j]: return False  
    else: return less(x[i+1:], x[j+1:])
```

First attempt



```
def less(x[i:], x[j:]):  
    if x[i] < x[j]: return True  
    if x[i] > x[j]: return False  
    else: return less(x[i+1:], x[j+1:])
```

First attempt

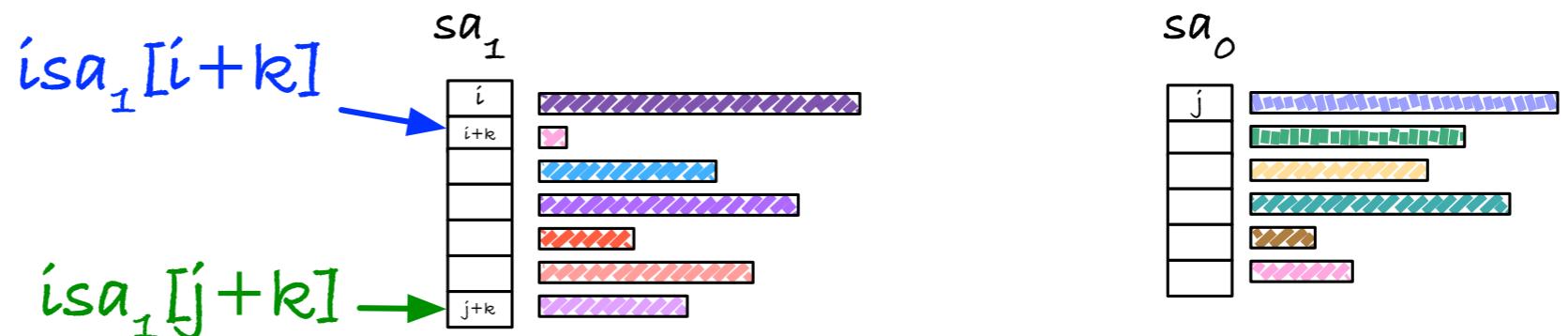


If we keep going, comparing identical letters, the comparison takes $O(n)$ time, so the merge takes $O(n^2)$.

Not cool!!!

First attempt

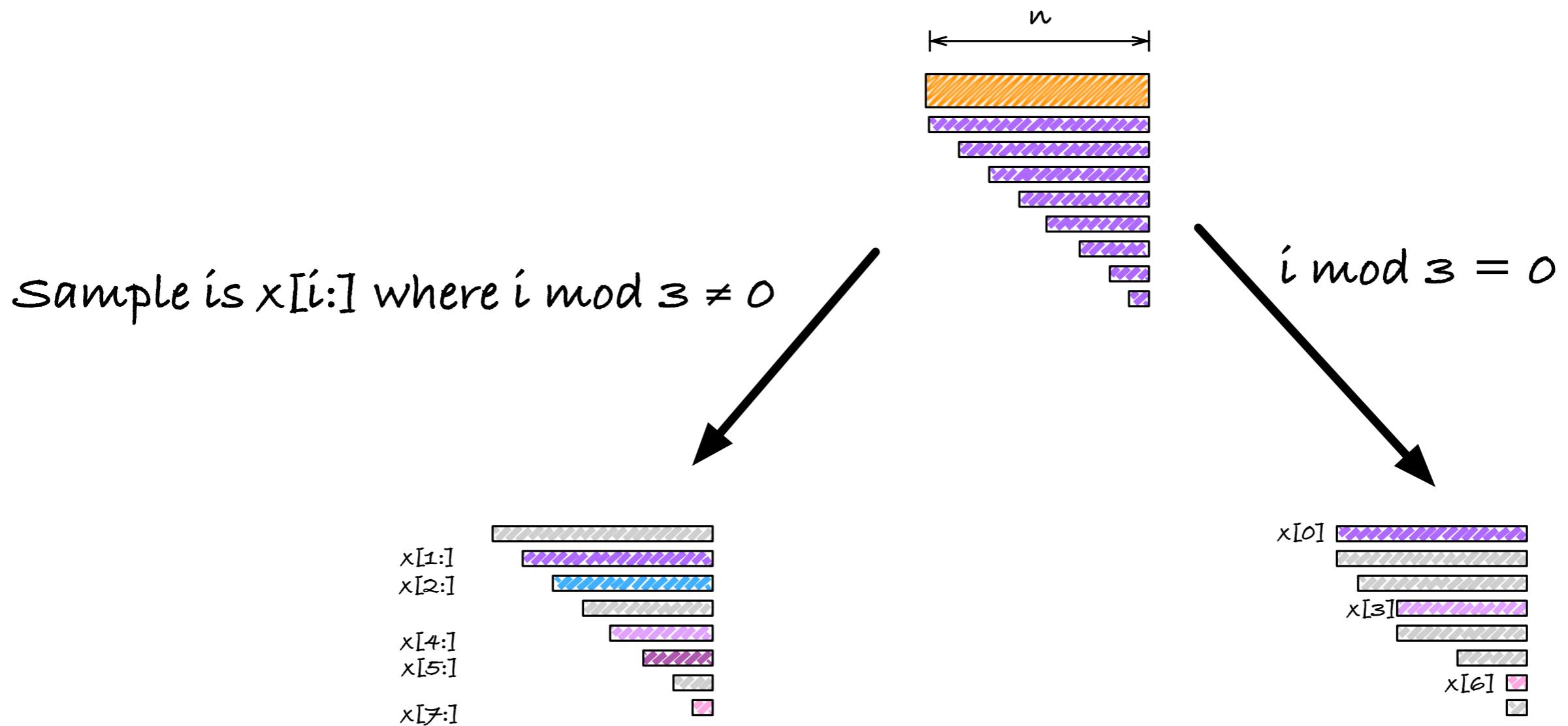
Wouldn't it be nice if we could get both indices in the same array?
Then we could determine their order in $O(1)$



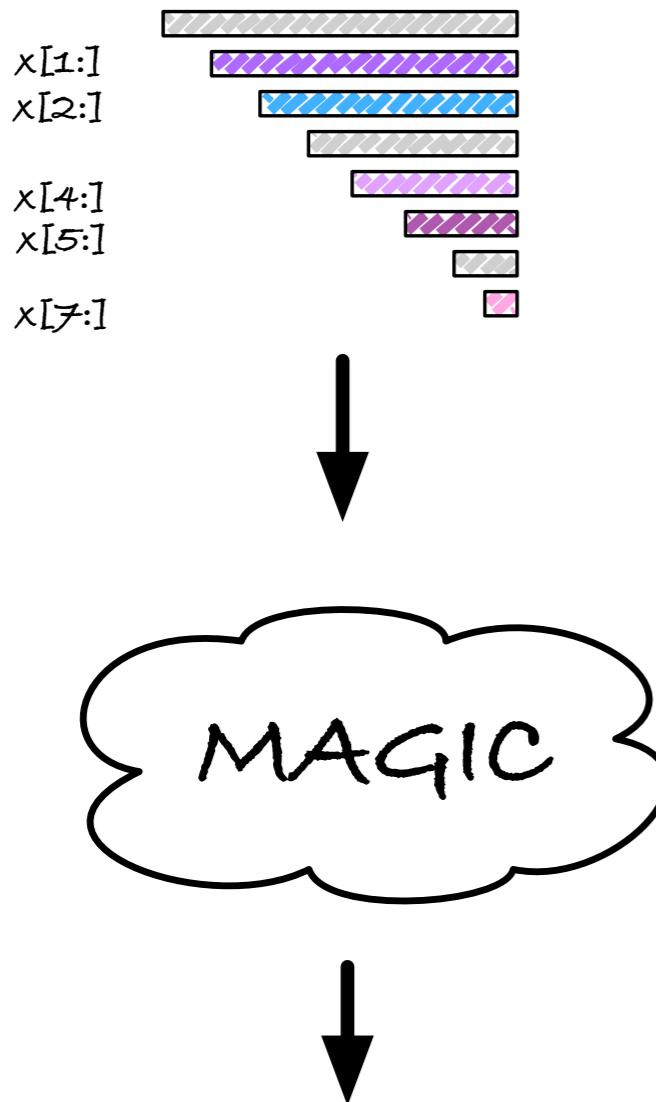
Well, we can't. The isa pointers will keep swapping array and never end up on the same side.

Not unless we make the sample and out-of-sample different sizes!

Skew algorithm

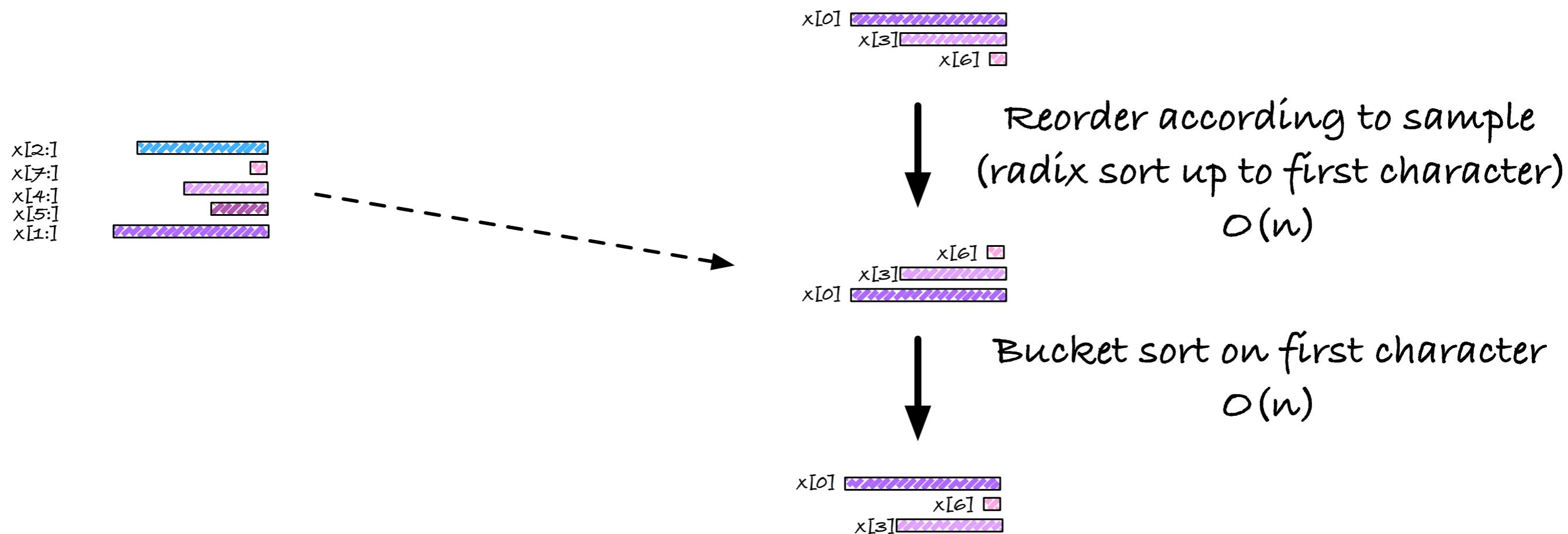


Skew algorithm



I will tell you how to do magic
when we are sure that we can merge...

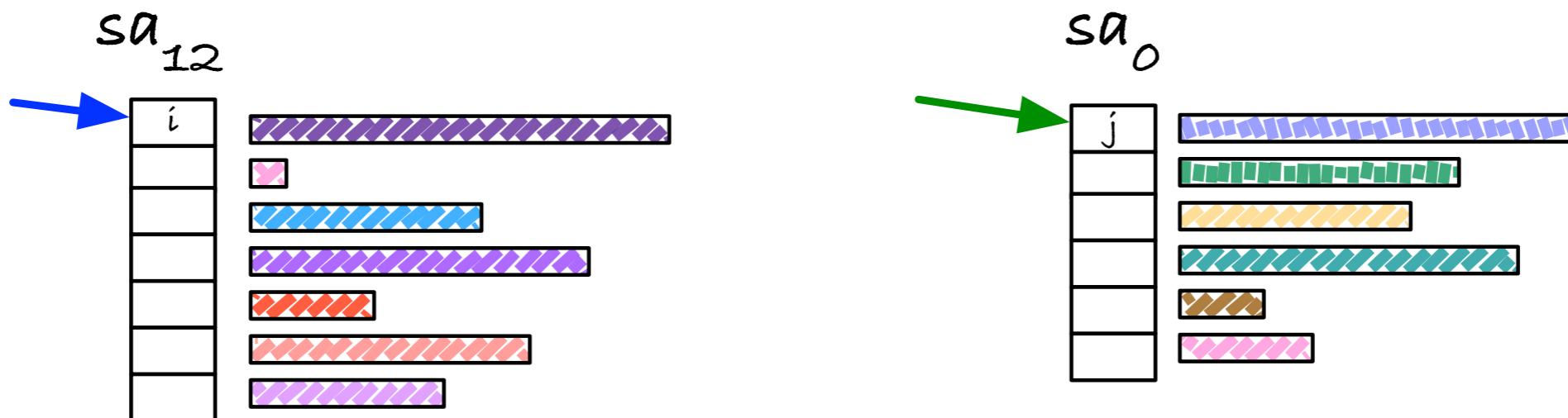
Skew algorithm



```
[i-1 for i in sa12 if i % 3 == 1]  
(plus special case if n % 3 == 0)
```

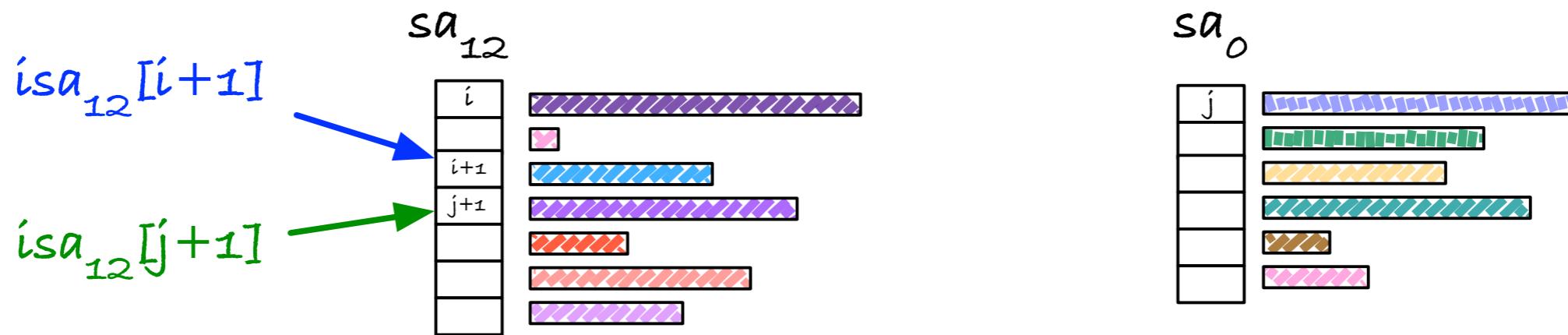
Skew algorithm

```
def less(x[i:],x[j:]):  
    if x[i] < x[j]: return True  
    if x[i] > x[j]: return False  
    else: return less(x[i+1:],x[j+1:])
```



Skew algorithm

```
def less(x[i:],x[j:]):  
    if x[i] < x[j]: return True  
    if x[i] > x[j]: return False  
    else: return less(x[i+1:],x[j+1:])
```

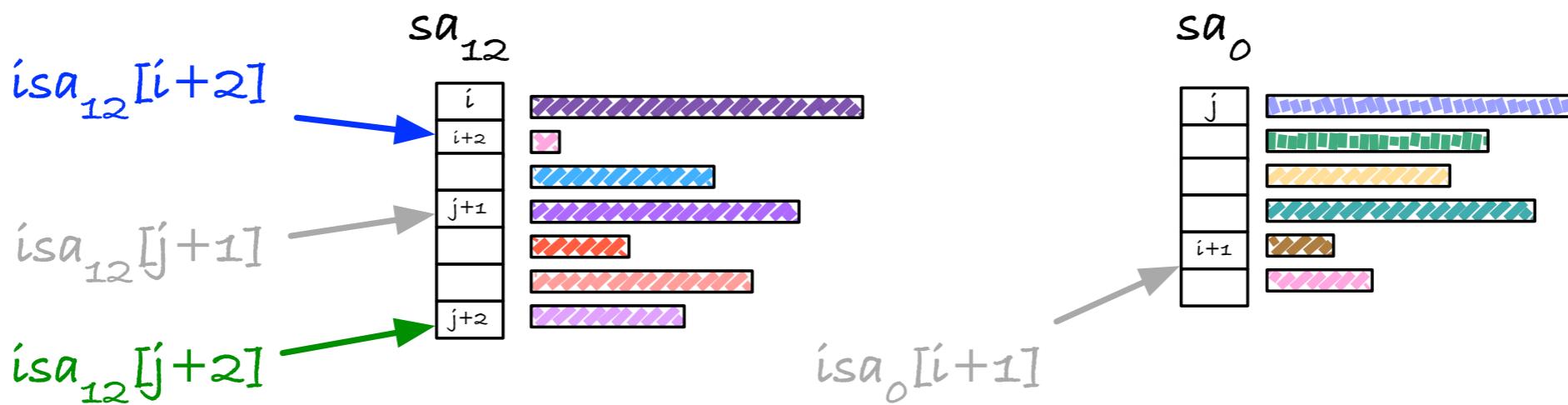


$$(j+1) \bmod 3 = 1$$

if $i \bmod 3 = 1$ then $(i+1) \bmod 3 = 2$

Skew algorithm

```
def less(x[i:],x[j:]):  
    if x[i] < x[j]: return True  
    if x[i] > x[j]: return False  
    else: return less(x[i+1:],x[j+1:])
```



$(j+1) \bmod 3 = 1$ and $(j+2) \bmod 3 = 2$

if $i \bmod 3 = 2$ then $(i+1) \bmod 3 = 0$

and $(i+2) \bmod 3 = 1$

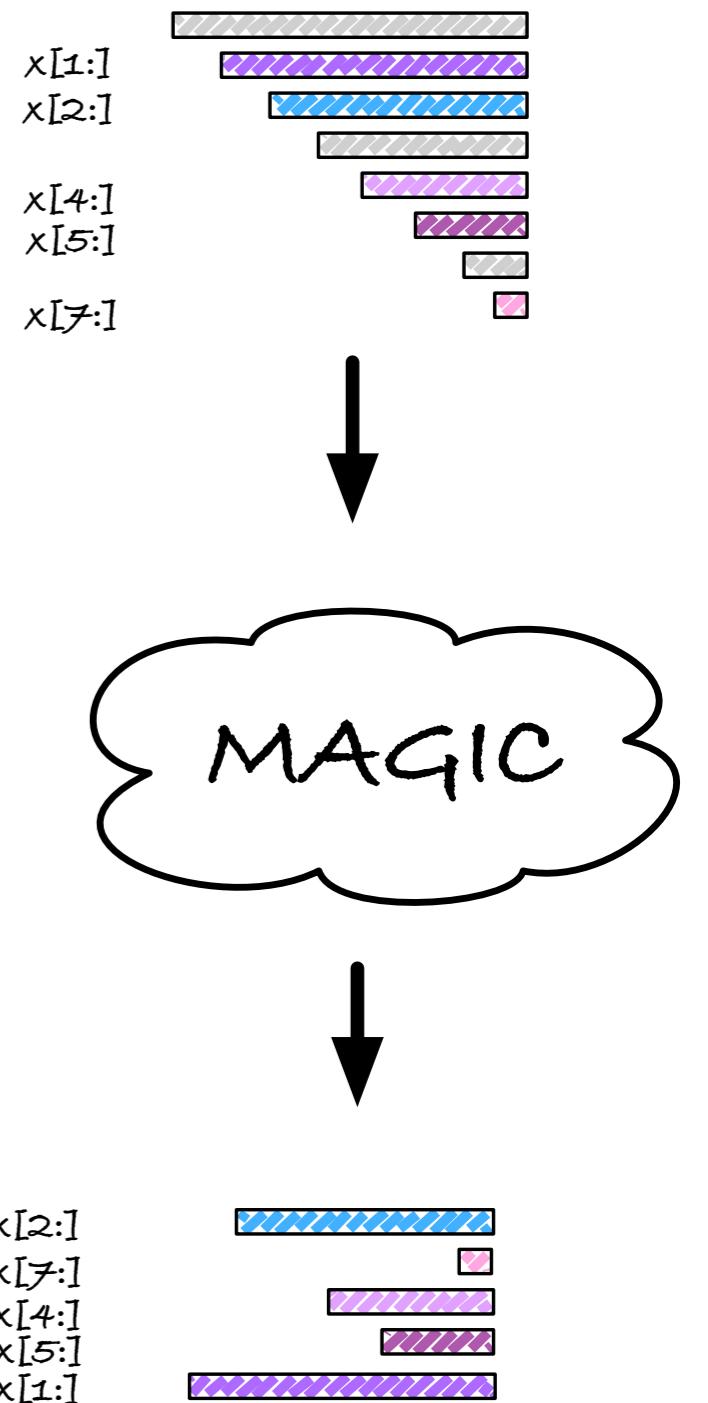
Skew algorithm

- In no more than two steps, you can determine $\text{less}(x[i:], x[j:])$, so comparison is in $O(1)$
- Therefore, the merge is in $O(n)$

MAGIC!!!

- **What do we need for the magic step?**

- We must create a new string, u , from the samples (the algorithm takes a string as input, so we can only recurse with a string)
- We must do that in $O(n)$ (or the recurrence doesn't work)
- The suffix array for u must give us the order of the samples



Constructing u

**This is going to look a bit weird,
but we will get to the good place
in the end...**

```
u = x[1:] # x[2:]
```

```
if x = "mississippi" then  
u = "ississippi#ssissippi"
```

u = x[2:]\#x[1:] works just as well.

(you need to change some indices later in the algorithm,
but which of x[1:] or x[2:] you put first doesn't matter)

Constructing u

This is going to look a bit weird,
but we will get to the good place
in the end...

The # character is the sentinel
(smaller than all other characters)

$u = x[1:] \# x[2:]$

It is a different symbol than \$ because it serves another purpose
(the “central” versus “terminal” sentinel)
but it *is* the same sentinel when you implement it.

In the skew algorithm you don’t need the terminal sentinel in the input
(although we use it implicitly when radix sorting)

So, we don’t need to add \$ to the end of u here.

Constructing u

**This is going to look a bit weird,
but we will get to the good place
in the end...**

```
u = x[1:] # x[2:]
```

```
if x = "mississippi" then  
u = "ississippi#ssissippi"
```

The length of u is not smaller than x!

$$|u| = 2n-2$$

Recurising on a *longer* string is folly!

Constructing u

**This is going to look a bit weird,
but we will get to the good place
in the end...**

```
u = x[1:] # x[2:]
```

```
if x = "mississippi" then  
u = "ississippi#ssissippi"
```

Group in triplets:

```
u = [iss] [iss] [ipp] [i$$] [#] [ssi] [ssi] [ppi]
```

Constructing u

**This is going to look a bit weird,
but we will get to the good place
in the end...**

If we can create an alphabet of triplets, and create u as a string with this alphabet, all in $O(n)$, then this u has length $\frac{2}{3}n$ and then we can recurse with $T(n)=T(\frac{2}{3}n)+O(n)$.

$u = x[1]$

Put a pin in this, but remember that we need to create a triplet alphabet.

if $x = \text{"mississippi"}$
 $u = \text{"ississippi"}$

Let's first figure out what we can do with u, though, to see if it is worthwhile to do it.

Group in triplets:

$u = [\text{iss}] [\text{iss}] [\text{ipp}] [\text{i\$\$}] [\#] [\text{ssi}] [\text{ssi}] [\text{ppi}]$

Constructing u

Consider the suffixes of u...

u:	x[1:]:	[íss] [íss] [ípp] [í\$\$\$] [#] [ssi] [ssi] [ppú]
u[0:]	x[1:]:	[íss] [íss] [ípp] [í\$\$\$] [#] [ssi] [ssi] [ppú]
u[1:]	x[4:]:	[íss] [ípp] [í\$\$\$] [#] [ssi] [ssi] [ppú]
u[2:]	x[7:]:	[ípp] [í\$\$\$] [#] [ssi] [ssi] [ppú]
u[3:]	x[10:]:	[í\$\$\$] [#] [ssi] [ssi] [ppú]
u[4:]		[#] [ssi] [ssi] [ppú]
u[5:]	x[2:]:	[ssi] [ssi] [ppú]
u[6:]	x[5:]:	[ssi] [ppú]
u[7:]	x[8:]:	[ppú]

Constructing u

Consider the suffixes of u...

$u:$		$[\acute{i}ss] \ [i\acute{ss}] \ [\acute{i}pp] \ [\acute{i}\$\$] \ [\#] \ [ssi] \ [ssi] \ [ppi]$
$u[0:]$	$x[1:]$:	$[\acute{i}ss] \ [i\acute{ss}] \ [\acute{i}pp] \ [\acute{i}\$\$] \ [\#] \ [ssi] \ [ssi] \ [ppi]$
$u[1:]$	$x[4:]$:	$[\acute{i}ss] \ [\acute{i}pp] \ [\acute{i}\$\$] \ [\#] \ [ssi] \ [ssi] \ [ppi]$
$u[2:]$	$x[7:]$:	$[\acute{i}pp] \ [\acute{i}\$\$] \ [\#] \ [ssi] \ [ssi] \ [ppi]$
$u[3:]$	$x[10:]$:	$[\acute{i}\$\$] \ [\#] \ [ssi] \ [ssi] \ [ppi]$
$u[4:]$		$[\#] \ [ssi] \ [ssi] \ [ppi]$
$u[5:]$	$x[2:]$:	$[ssi] \ [ssi] \ [ppi]$

When you compare $u[i:]$ with $u[j:]$ for some i,j , you never compare past the # sentinel (they always differ there, with # being the smallest letter).

The order of suffixes $u[i:]$ does not depend on the strings past the sentinel. The $x[i:] \# x[2:]$ ($i \bmod 3 = 1$) strings are in effect just the $x[i:]$ strings.

Constructing u

However! u is over a different alphabet. All the triplets are really different symbols (in practise, different integers)

$$[xyz] = \xi$$

[ssi] [ssi] [ppi]

[ssi] [ssi] [ppi]

$$u[2:] = \xi\omega = x[2:] = [ipp] [i\$\$] [\i\$\$] [\#] [ssi] [ssi] [ppi]$$

$$u[2:] \quad x[7:]:$$

[ipp] [*\\$*\\$] [\#] [ssi] [ssi] [ppi]

$$u[3:]$$

If the order of u suffixes should match the order of the

corresponding x suffixes, then the alphabet must be order-preserving,

$$u[5:]$$

i.e., $[xyz] < [abc]$ iff “xyz” < “abc”.

$$u[6:] \quad x[5:]:$$

[ssi] [ppi]

$$u[7:] \quad x[8:]:$$

[ppi]

Creating the alphabet

- We have to create an alphabet of triplets
 - I.e. each triplet must map to a unique integer
- The mapping must be order-preserving, $[xyz] < [abc]$ iff “xyz” < “abc”
 - (This is usually the tricky part, since it is easy to build a table of triplets, but sorting them might not be)
- We only have $O(n)$ time to build the alphabet in, or the running time breaks

Creating the alphabet

Take the $i \bmod 3 \neq 0$ suffixes and radix sort them on the first three letters. Then you have the sorted triplets.

$x[1:]$:	iss	iss	ipp	i		$x[10:]$:	i\$	\$		
$x[4:]$:	iss	ipp	i			$x[2:]$:	ssi	ssi	ppi	
$x[7:]$:	ipp	i				$x[5:]$:	ssi	ppi		
$x[10:]$:	i\$	\$				$x[8:]$:	pp	i		
$x[2:]$:	ssi	ssi	ppi			$x[7:]$:	ipp	i		
$x[5:]$:	ssi	ppi				$x[1:]$:	iss	iss	ipp	i
$x[8:]$:	ppi					$x[4:]$:	iss	iss	ipp	i

Runs in $O(n + \sigma)$

$x[10:]$:	i\$	\$			$x[10:]$:	i\$	\$		
$x[2:]$:	ssi	ssi	ppi		$x[8:]$:	pp	i		
$x[5:]$:	ssi	ppi			$x[7:]$:	ipp	i		
$x[8:]$:	ppi				$x[2:]$:	ssi	ssi	ppi	
$x[7:]$:	ipp	i			$x[5:]$:	ssi	ppi		
$x[1:]$:	iss	iss	ipp	i	$x[1:]$:	iss	iss	ipp	i
$x[4:]$:	iss	ipp	i		$x[4:]$:	iss	ipp	i	

But we need $O(n)$... is $O(n + \sigma)$ in $O(n)$?

$x[10:]$:	i\$	\$			$x[10:]$:	i\$	\$		
$x[8:]$:	ppi				$x[7:]$:	ipp	i		
$x[7:]$:	ipp	i			$x[1:]$:	iss	iss	ipp	i
$x[2:]$:	ssi	ssi	ppi		$x[4:]$:	iss	ipp	i	
$x[5:]$:	ssi	ppi			$x[8:]$:	pp	i		
$x[1:]$:	iss	iss	ipp	i	$x[2:]$:	ssi	ssi	ppi	
$x[4:]$:	iss	ipp	i		$x[5:]$:	ssi	ppi		

Worries about the alphabet size?

- In this algorithm (and others) we create new alphabets. They are not constant sized! So, we need to consider the alphabet size as part of the time analysis.
- However, we can't have more than $O(n)$ triplets if we get them from $\frac{2}{3}n$ different strings. So the alphabet size, σ , is in $O(n)$
- This is a far cry from $O(1)$, but we only see σ in the bucket sorts in this algorithm, which runs in $O(n + \sigma)$, which is then in $O(n)$.
- We are okay in this algorithm, even if the alphabet size isn't constant.

Creating the alphabet

{“#”: 0}

x[10:]: i\$\$
x[7:]: ipp i
x[1:]: iss iss ipp i
x[4:]: iss ipp i
x[8:]: ppi
x[2:]: ssi ssi ppi
x[5:]: ssi ppi

Creating the alphabet

```
{ "#": 0,  
  "i$$": 1,  
}
```



```
x[10:]: i$$  
x[7:]: ipp i  
x[1:]: iss iss ipp i  
x[4:]: iss ipp i  
x[8:]: ppi  
x[2:]: ssi ssi ppi  
x[5:]: ssi ppi
```

Creating the alphabet

```
{ "#": 0,  
  "i$$": 1,  
  "ipp": 2  
}
```



$x[10:]$:	i\$\$
$x[7:]$:	ipp i
$x[1:]$:	iss iss ipp i
$x[4:]$:	iss ipp i
$x[8:]$:	ppi
$x[2:]$:	ssi ssi ppi
$x[5:]$:	ssi ppi

Creating the alphabet

```
{ "#": 0,  
  "i$$": 1,  
  "ipp": 2,  
  "iss": 3  
}
```

→

$x[10:]$:	i\$\$
$x[7:]$:	ipp i
$x[1:]$:	iss iss ipp i
$x[4:]$:	iss ipp i
$x[8:]$:	ppi
$x[2:]$:	ssi ssi ppi
$x[5:]$:	ssi ppi

Creating the alphabet

```
{ "#": 0,  
  "i$$": 1,  
  "ipp": 2,  
  "iss": 3  
}
```

$x[10:]$: i\$\$
 $x[7:]$: ipp i
 $x[1:]$: iss iss ipp i
 $x[4:]$: iss ipp i
 $x[8:]$: ppi
 $x[2:]$: ssi ssi ppi
 $x[5:]$: ssi ppi



Creating the alphabet

```
{ "#": 0,  
  "i$$": 1,  
  "ipp": 2,  
  "iss": 3,  
  "ppi": 4  
}
```

$x[10:]$: i\$\$
 $x[7:]$: ipp i
 $x[1:]$: iss iss ipp i
 $x[4:]$: iss ipp i
 $x[8:]$: ppi
 $x[2:]$: ssi ssi ppi
 $x[5:]$: ssi ppi



Creating the alphabet

```
{ "#": 0,  
  "i$$": 1,  
  "ipp": 2,  
  "iss": 3,  
  "ppi": 4,  
  "ssi": 5  
}
```

$x[10:]$: i\$\$
 $x[7:]$: ipp i
 $x[1:]$: iss iss ipp i
 $x[4:]$: iss ipp i
 $x[8:]$: ppi
 $x[2:]$: ssi ssi ppi
 $x[5:]$: ssi ppi



Creating the alphabet

```
{ "#": 0,  
  "i$$": 1,  
  "ipp": 2,  
  "iss": 3,  
  "ppi": 4,  
  "ssi": 5  
}
```

$x[10:]$: i\$\$
 $x[7:]$: ipp i
 $x[1:]$: iss iss ipp i
 $x[4:]$: iss ipp i
 $x[8:]$: ppi
 $x[2:]$: ssi ssi ppi
 $x[5:]$: ssi ppi



Creating the alphabet

- Because we sort the triplets with radix sort, we can preserve the order in $O(n)$.
- If you just created the table and sorted the keys, then you would pay for the sorting, typically $O(n \log n)$
 - $O(n \log n)$ comparisons with comparisons of triplets taking constant time

You have the alphabet, now what?

- This is an excellent place to consider how we stop the recursion!

You have the alphabet, now what?

- This is an excellent place to consider how we stop the recursion!
- You did wonder, didn't you? When we use recursion, we need a base case where we stop (and we haven't talked about that at all)

You have the alphabet, now what?

- This is an excellent place to consider how we stop the recursion!
 - You did wonder, didn't you? When we use recursion we need a base case where we stop (and we haven't talked about that at all)
-
- ```
graph LR; A[skew(x)] --> B[skew(u)]; B --> C[skew(u)]; C --> D[skew(u)]; D --> E[infinity]
```
- The diagram illustrates a recursive call. It starts with a yellow box labeled "skew(x)". An arrow points from "skew(x)" to a yellow box labeled "skew(u)". From this "skew(u)" box, two arrows point to two separate yellow boxes, both labeled "skew(u)". Finally, an arrow points from one of these "skew(u)" boxes to a red box containing the symbol "
- $\infty$
- ".

# Base case

- The purpose of building  $u$  is to sort suffixes  $i \bmod 3 \neq 0$
- If all leading triplets are unique, they are already sorted
  - We can determine the order of strings from the prefix until they differ
  - If they differ in the first three letters, sorting on the first three letter gives us the correct order
  - Then we don't need to recurse and can stop here!

# Constructing u

- Easy way to build u is to put the triplets and their numbers in a table, then straight-forward scan through  $x[1+i:1+i+3]$  and then  $x[2+i:2+i+3]$  in jumps of 3

```
[*(alpha[triplet(x, i)] for i in range(1, len(x), 3)),
SENTINEL,
*(alpha[triplet(x, i)] for i in range(2, len(x), 3))]
```

# Recurse, and then what?

- You know how to map between indices in  $x[i:]$   $i \bmod 3 \neq 0$  and  $u[j:]$  because you created  $u$  in a specific way.
- Use that to map  $u[j:]$  back to  $x[i:]$

|        |            |                                                                                                                     |
|--------|------------|---------------------------------------------------------------------------------------------------------------------|
| $sa_u$ | $u[4:]$    | $\text{[#] } [\text{ssí}] \ [ \text{ssí} ] \ [\text{ppí}]$                                                          |
|        | $u[3:]$    | $[\text{í\$\$}] \ [\#] \ [\text{ssí}] \ [ \text{ssí} ] \ [\text{ppí}]$                                              |
|        | $u[2:]$    | $[\text{ípp}] \ [\text{í\$\$}] \ [\#] \ [\text{ssí}] \ [ \text{ssí} ] \ [\text{ppí}]$                               |
|        | $u[1:]$    | $[\text{íss}] \ [\text{ípp}] \ [\text{í\$\$}] \ [\#] \ [\text{ssí}] \ [ \text{ssí} ] \ [\text{ppí}]$                |
|        | $u[0:]$    | $[\text{íss}] \ [\text{íss}] \ [\text{ípp}] \ [\text{í\$\$}] \ [\#] \ [\text{ssí}] \ [ \text{ssí} ] \ [\text{ppí}]$ |
|        | $u[7:]$    | $[\text{ppí}]$                                                                                                      |
|        | $u[6:]$    | $[\text{ssí}] \ [\text{ppí}]$                                                                                       |
|        | $u[5:]$    | $[\text{ssí}] \ [\text{ssí}] \ [\text{ppí}]$                                                                        |
|        |            | $0 \quad 1 \quad 2 \quad 3 \quad 4 \quad 5 \quad 6 \quad 7$                                                         |
|        |            | $u = x[1...]\times[4...]\times[7...]\times[10...]\#x[2...]\times[5...]\times[8...]$                                 |
|        | $x[10:]$ : | $\text{i}$                                                                                                          |
|        | $x[7:]$ :  | $\text{íppí}$                                                                                                       |
|        | $x[4:]$ :  | $\text{íssíppí}$                                                                                                    |
|        | $x[1:]$ :  | $\text{íssíssíppí}$                                                                                                 |
|        | $x[8:]$ :  | $\text{ppí}$                                                                                                        |
|        | $x[5:]$ :  | $\text{ssíppí}$                                                                                                     |
|        | $x[2:]$ :  | $\text{ssíssíppí}$                                                                                                  |

# Recurse, and then what?

- You know how to map between indices in  $x[i:]$   $i \bmod 3 \neq 0$  and  $u[j:]$  because you created  $u$  in a specific way.

$sa_u$

|         |                                                                                                                 |
|---------|-----------------------------------------------------------------------------------------------------------------|
| $u[4:]$ | $\text{[#] } [\text{ssí}] \ [ \text{ssí}] \ [\text{ppí}]$                                                       |
| $u[3:]$ | $[\text{###}] \ [\#] \ [\text{ssí}] \ [\text{ssí}] \ [\text{ppí}]$                                              |
| $u[2:]$ | $[\text{ppí}] \ [\text{###}] \ [\#] \ [\text{ssí}] \ [\text{ssí}] \ [\text{ppí}]$                               |
| $u[1:]$ | $[\text{ssí}] \ [\text{ppí}] \ [\text{###}] \ [\#] \ [\text{ssí}] \ [\text{ssí}] \ [\text{ppí}]$                |
| $u[0:]$ | $[\text{ssí}] \ [\text{ssí}] \ [\text{ppí}] \ [\text{###}] \ [\#] \ [\text{ssí}] \ [\text{ssí}] \ [\text{ppí}]$ |
| $u[7:]$ | $[\text{ppí}]$                                                                                                  |
| $u[6:]$ | $[\text{ssí}] \ [\text{ppí}]$                                                                                   |
| $u[5:]$ | $[\text{ssí}] \ [\text{ssí}] \ [\text{ppí}]$                                                                    |

$$u = x[1\ldots]x[4\ldots]x[7\ldots]x[10\ldots] \# x[2\ldots]x[5\ldots]x[8\ldots]$$

- Use that to map  $u[j:]$  back to  $x[i:]$

$x[10:]$ :  $i$

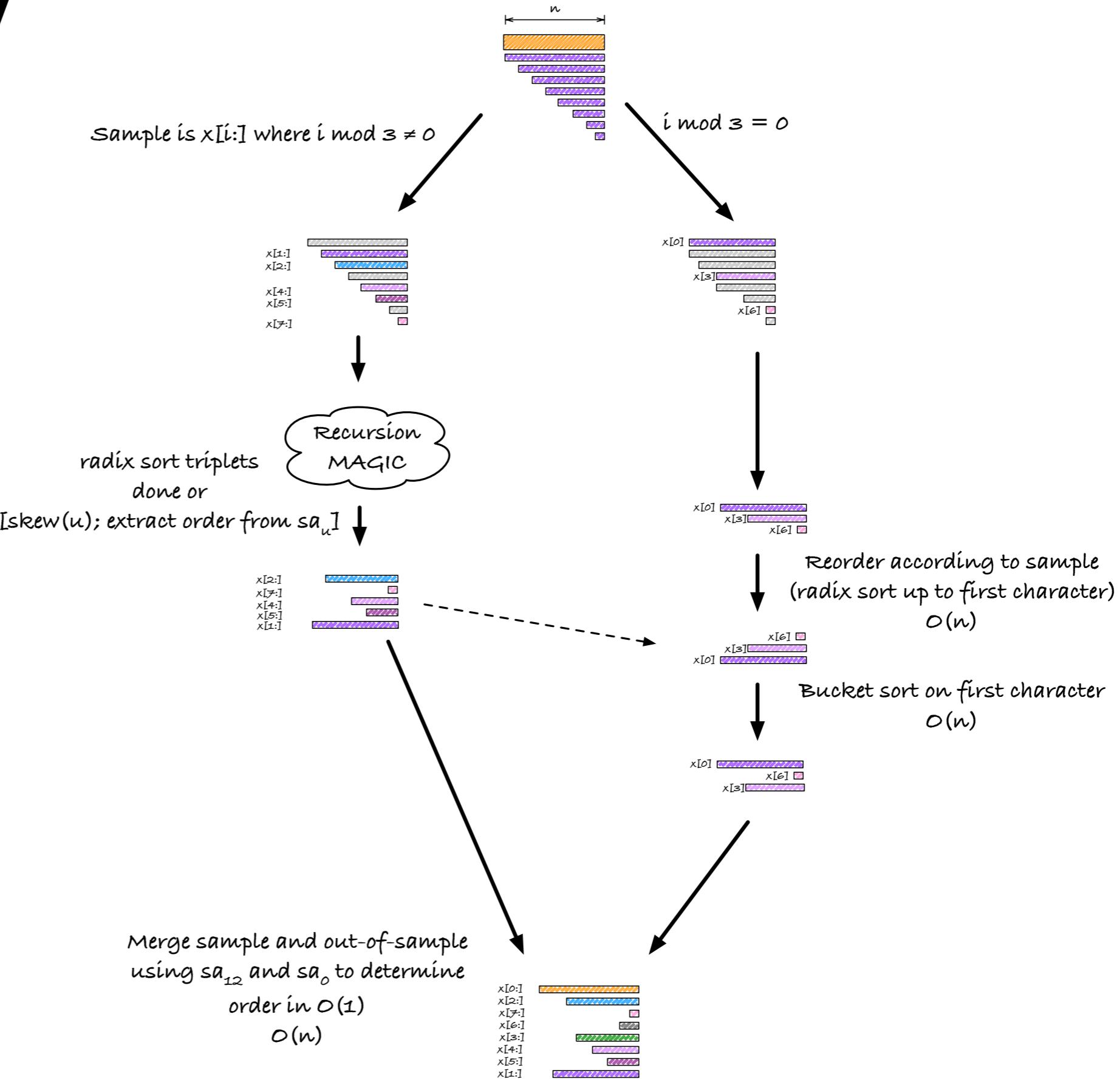
$x[7:]$ :  $\text{íppí}$

The exact mapping from  $u[j:]$  to  $x[i:]$  depends on how you construct  $u$  ( $u = x[1:] \# x[2:]$  or  $x[2:] \# x[1:]$ ) or whether you include the central sentinel (you can avoid it if you start with a terminal \$), but the mapping is always just a bit of arithmetic.

# Summary of magic

- Radix sort the leading triplet of the suffixes  $i \bmod 3 \neq 0$ 
  - If the triplets are unique, we have sorted the suffixes
  - Otherwise,
    - construct  $u$  from triplets in  $x[1:] \# x[2:]$
    - call  $\text{skew}(u)$
    - map the indices in  $\text{sa}_u$  back to the suffixes  $i \bmod 3 \neq 0$  (this depends on how we constructed  $u$ , but is straightforward)

# Summary of skew



# Final words on sentinels

- The central sentinel # when we construct  $u$  separates the  $i \bmod 3 = 1$  from the  $i \bmod 3 = 2$  strings. It must be there, or the algorithm fails.
- However, if the input,  $x$ , has a terminal sentinel, then the triplets in the alphabet we generate from  $x$  will contain a symbol that functions as #.
  - If we create  $u = x[1:]x[2:]$  or  $u = x[2:]x[1:]$  (just concatenate the two classes without the # sentinel), the center triplet will be unique and function as #.
  - This is a tedious proof, essentially just case analysis, and I won't show it.

*That's all folks!*