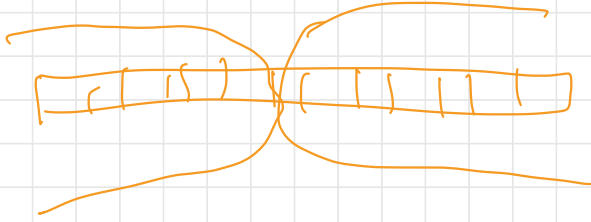
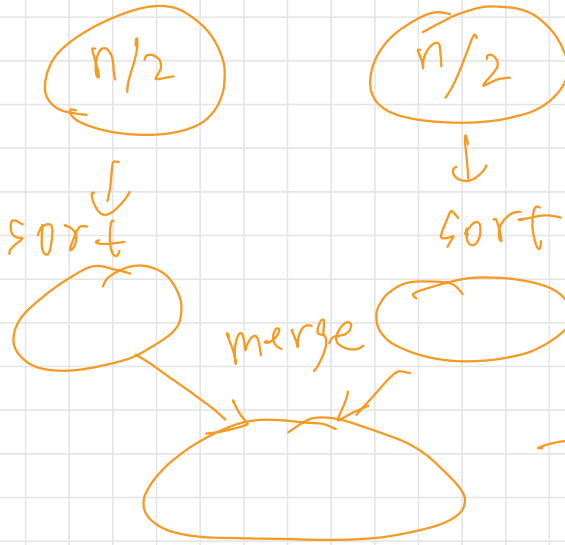


# Review Merge Sort



Divide.



$$T(n)$$

$$= 2 \cdot T(n/2)$$

$$+ n$$



$$T(n) = O(n \cdot \log n)$$

# Skew Algorithm For Suffix Sorting

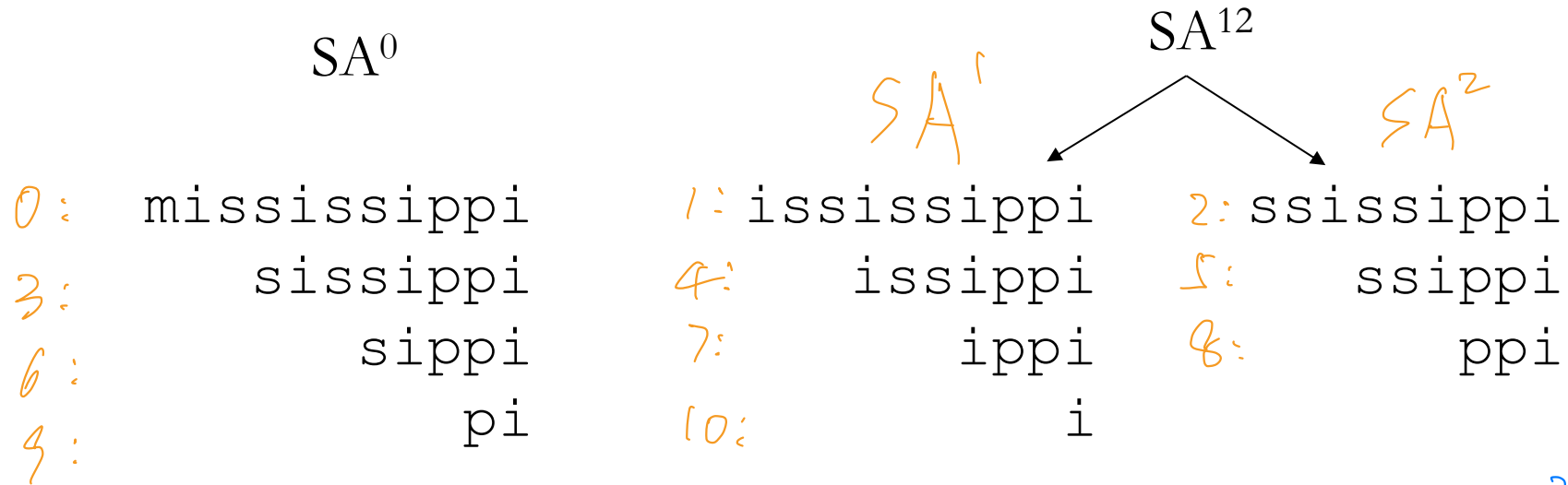
---

- Let  $S_0, S_1, S_2, \dots, S_{n-1}$  be all the  $n$  suffixes.  $S_i$  starts at  $i$ -th position.
- Skew algorithm uses divide and conquer. But it divides the problem into unequally sized parts.
- Two sets  $SA^0 = \{S_i : i = 0 \pmod{3}\}$  and  $SA^{12} = \{S_i : i = 1 \text{ or } 2 \pmod{3}\}$ .

# Skew Algorithm Example

- Example: mississippi

0 1 2 3 4 5 6 7 8 9 ⑩

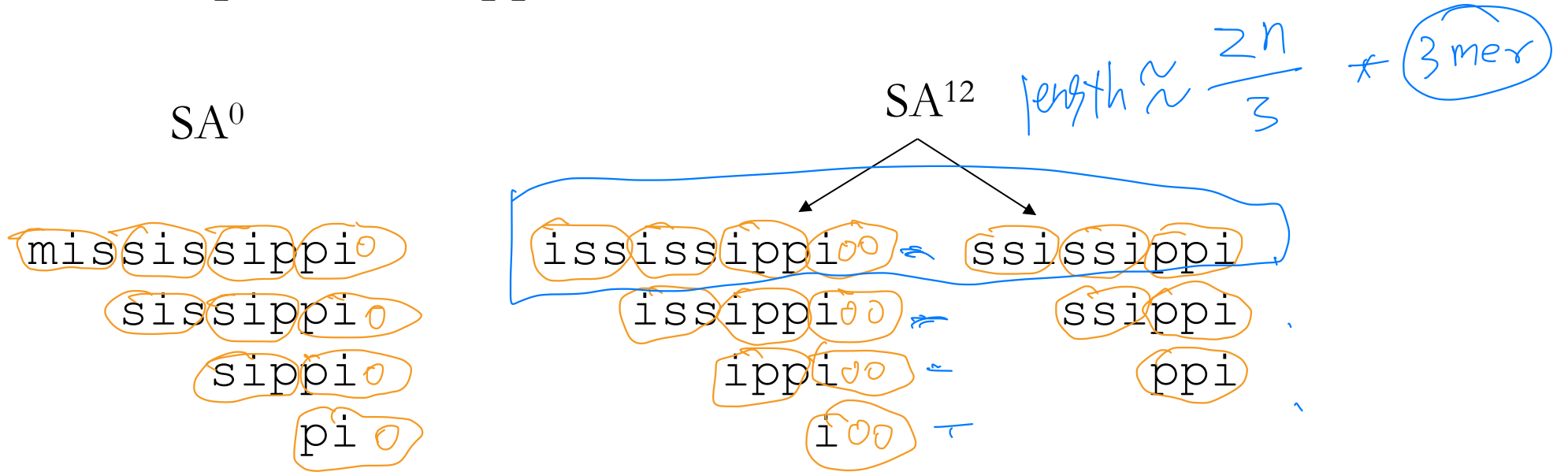


- ✓ ① Sort  $SA^{12}$  recursively
- ✓ ② Sort  $SA^0$  in linear time
- ✓ ③ merge in linear time

$$T(n) = T\left(\frac{2}{3}n\right) + n + n = O(n).$$

# Skew Algorithm Example

- Example: mississippi



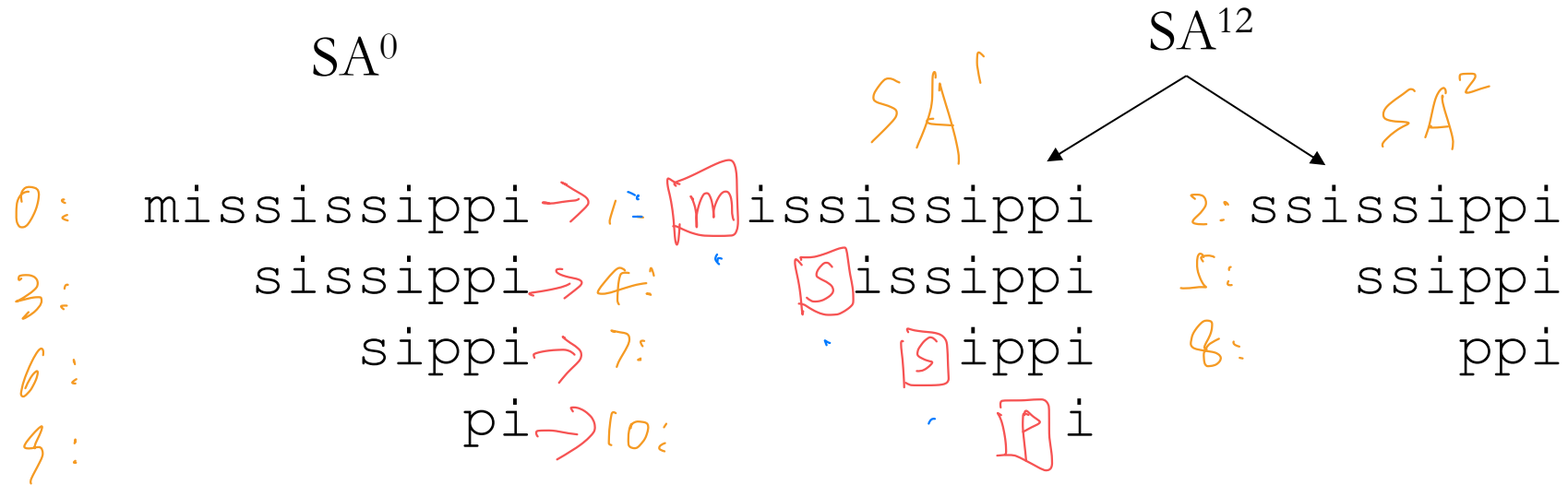
Renaming:  
3-mer → integer

Now we can do recursion,  
except 3-mer.

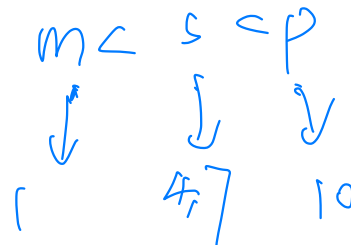
# Skew Algorithm Example

- Example: mississippi

0 1 2 3 4 5 6 7 8 9 10



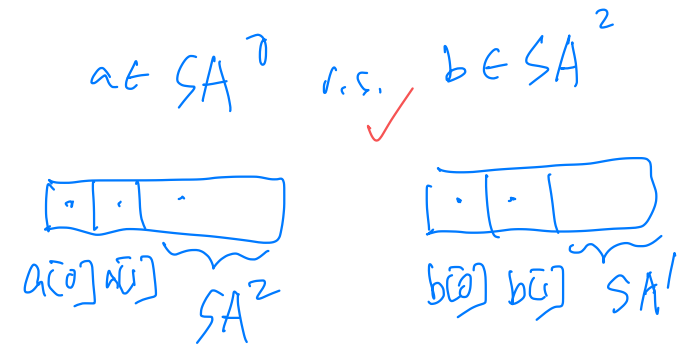
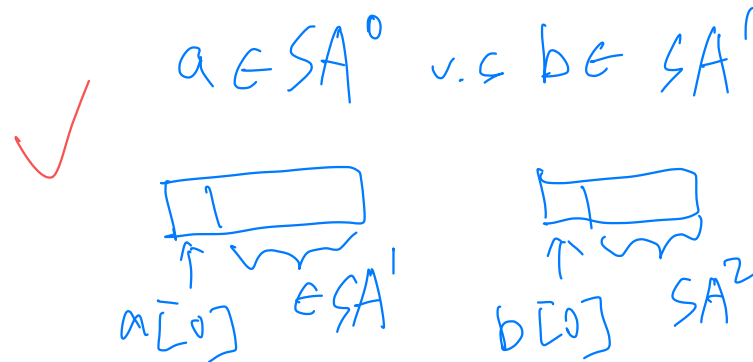
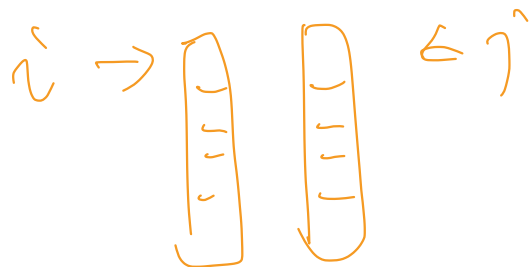
one-pass in radix sorting



# Skew Algorithm Example

- Example: mississippi  
0 1 2 3 4 5 6 7 8 9 10

	$SA^0$		$SA^{12}$		
		$SA^1$		$SA^2$	
0:	mississippi	1:	ississippi	2:	ssissippi
3:	sissippi	4:	issippi	5:	ssippi
6:	sippi	7:	ippi	8:	ppi
9:	pi	10:	i		



# Skew Algorithm For Suffix Sorting

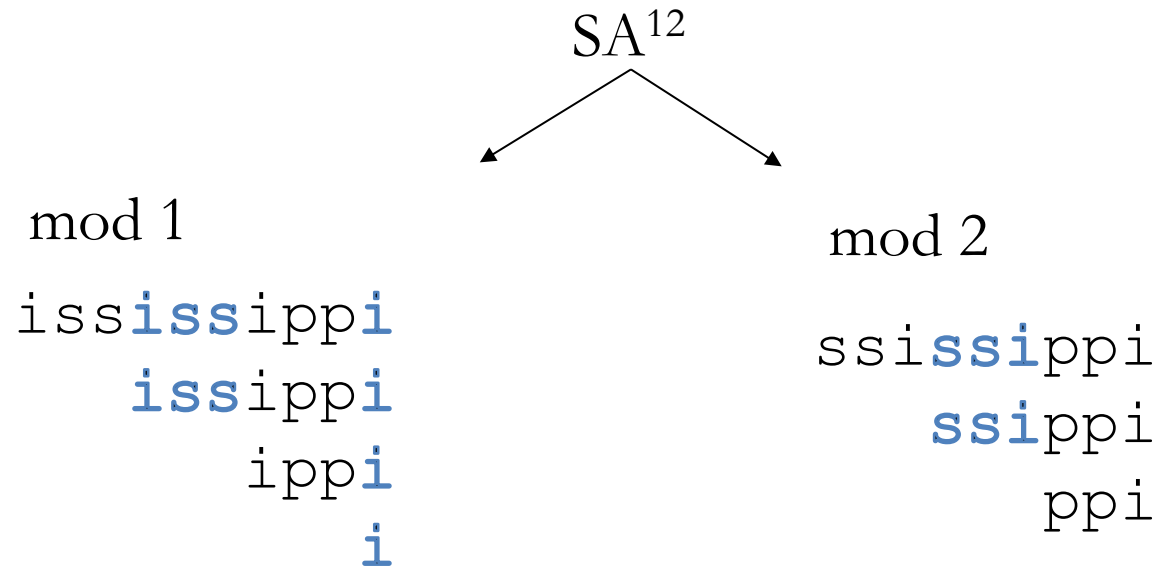
---

- 1. Sort  $SA^{12}$  recursively.
- 2. Sort  $SA^0$  in linear time.
- 3. Merge sort  $SA^0$  and  $SA^{12}$  in linear time.

# How to sort $SA^{12}$ recursively

---

mississippi



- We need to know the order of these suffixes.
- In order to solve it recursively, we need to reduce the problem to a smaller suffix sorting problem.



# Reduction to a smaller suffix sorting problem

---

SA<sup>12</sup>

mod 1

```
ississippi00
  issippi00
    ippip00
      i00
```

mod 2

```
ssissippi
  ssiippi
    ssiippi
      ppi
```

**i s s i s s i p p i 0 0**

**s s i s s i p p i**

- Pad 0 to make their length multiple of 3. Then treat each string as a string of “triplets”. Each subset is the suffixes of the “triplet string”.
- We connect the two “triplet strings” together to make a longer string. We put the one with padding at the left.

**i s s i s s i p p i 0 0 s s i s s i p p i**

# Reduction

---

SA<sup>12</sup>

mod 1

```

ississippi00
  issippi00
    ippi00
      i00
  
```

mod 2

```

ssissippi
  ssippi
    ppi
  
```

- Now check all the suffixes of the concatenated triplet string. Their relative order can be used to build the relative order of SA<sup>12</sup> easily.
- We are almost there, except that keeping tripling the size (number of bytes) of the “character” is a problem.

```

ississippi00ssissippi
  issippi00ssissippi
    ippi00ssissippi
      i00ssissippi
        ssissippi
          ssippi
            ppi
  
```

# Renaming

---

- We solve the unlimited expansion problem by a trick called renaming. It maps each unique triplet to a single unique integer.
- To rename, we first sort the triplets, and then assign integer values sequentially to unique triplets. Sorting triplets can be done in linear time by radix sort.
- This ensures
  - The max value is always bounded by the length of array.
  - The suffix order is unchanged.

i00 -> 0  
ipp -> 1  
iss -> 2  
ppi -> 3  
ssi -> 4

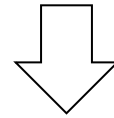
ississippi00ssiissippi  
↓  
2 2 1 0 4 4 3

# Renaming Example

---

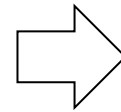
i00 -> 0  
ipp -> 1  
iss -> 2  
ppi -> 3  
ssi -> 4

ississippi00ssissippi



2 2 1 0 4 4 3

ississippi00ssissippi  
  issippi00ssissippi  
    ippi00ssissippi  
      i00ssissippi  
      ssiissippi  
      ssippi  
      ppi



2 2 1 0 4 4 3  
2 1 0 4 4 3  
1 0 4 4 3  
0 4 4 3  
4 4 3  
4 3  
3

# Recursion

---

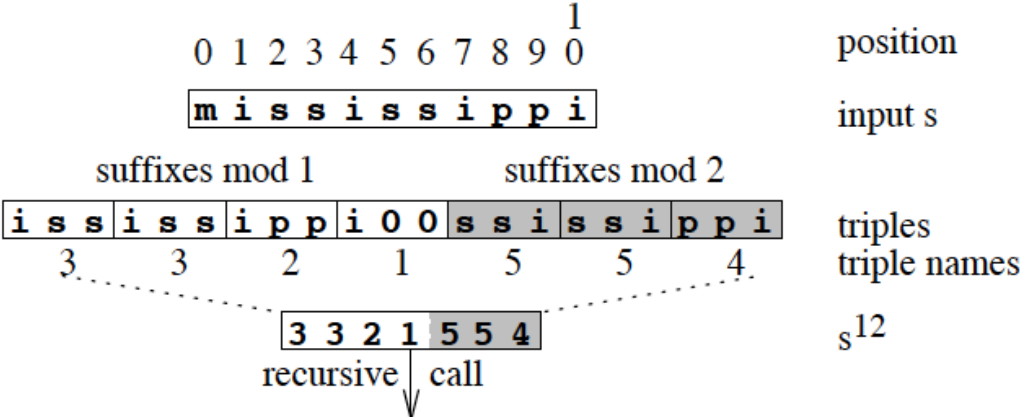
- After renaming, we just suffix sort the new integer string, which has length approximately  $2n/3$ . This can be done by recursion.
- The time complexity of renaming is dominated by sorting the triplets. This can be solved in linear time with radix sort.

# Radix Sort

---

- **Radix Sort:** Multiple passes. Each pass stable sorts according to one digit. From the least to the most significant digit.
- original: its, iss, ipp, abc, att
- pass1: abc, ipp, its, iss, att
- pass2: abc, ipp, iss, its, att
- pass3: abc, att, ipp, iss, its
  
- Radix sorting requires  $O(k)$  space, where  $k$  is the size of the alphabet.
- Each pass takes linear time. And only 3 passes needed in our case. So it is linear time.

# Recap Sort $S^{12}$ recursively



1. Padding and concatenation to get string of triplets.
2. Radix sort the triplets to get an ID (name) of each triple.
3. Recursion to get the suffix order on the string of IDs.

# Skew Algorithm For Suffix Sorting

---

- We assume  $SA^{12}$  is sorted already, and learn the other two steps first.
- 1. Sort  $SA^{12}$  recursively.
- 2. Sort  $SA^0$  in linear time.
- 3. Merge sort  $SA^0$  and  $SA^{12}$  in linear time.



# Sort $S^0$ in linear time

---

- $S_i = s[i] S_{i+1}$ .
- For all  $S_i$  in  $SA^0$ ,  $S_{i+1}$  has been sorted already. Use  $s[i]$  to do another pass of radix sorting will give us the right order of  $SA^0$ . This takes linear time.

Sorted  $SA_{12}$

	0	1	2	3	4	5	6	7	8	9	0 <sup>1</sup>
	<u>m i s s i s s i p p i</u>										

10: i  
4: issippi  
1: ississippi  
7: ippi  
8: ppi  
5: ssippi  
2: ssissippi

To sort  $SA_0$

0: **m**issippi  
3: **s**issippi  
6: **s**ippi  
9: **p**i

# Skew Algorithm For Suffix Sorting

---

- 1. Sort  $SA^{12}$  recursively.
- 2. Sort  $SA^0$  in linear time.
- 3. Merge sort  $SA^0$  and  $SA^{12}$  in linear time.

# Merge

---

Sorted SA12

10: i  
4: issippi  
1: ississippi  
7: ippi  
8: ppi  
5: ssippi  
2: ssissippi

Sorted SA0

0: mississippi  
9: pi  
6: sippi  
3: sissippi

- Would be a simple merge if comparison of two takes constant time.
- Trouble is when two suffices share a long prefix, which takes more than constant time to compare.

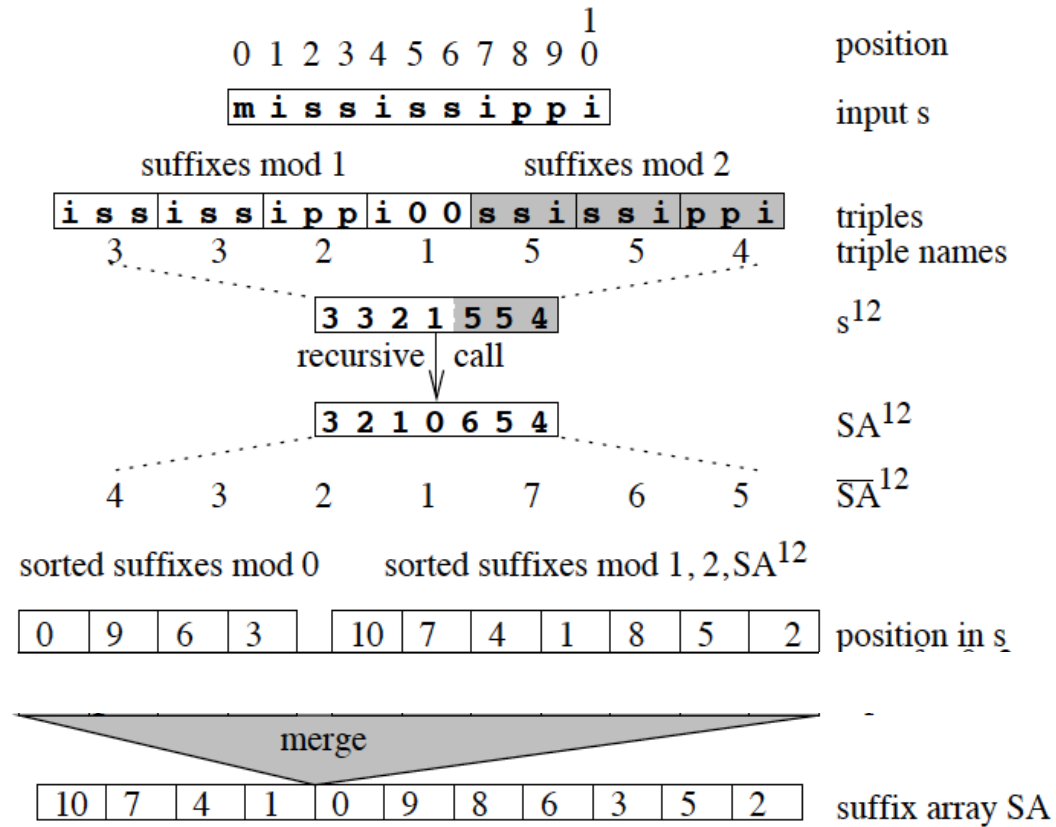
E.g. what if  $S_5 = \text{aaaa}\dots$  and  $S_6 = \text{aaaa}\dots$

## Merge $S^0$ and $S^{12}$

---

- Merging only requires to compare a suffix  $S_j$  with  $j \bmod 3 = 0$  with a suffix  $S_i$  with  $i \bmod 3 \neq 0$ . :
- Case 1: If  $i \bmod 3 = 1$ , we write  $S_i$  as  $(s[i], S_{i+1})$  and  $S_j$  as  $(s[j], S_{j+1})$ .
  - Since  $(i + 1) \bmod 3 = 2$  and  $(j + 1) \bmod 3 = 1$ , the relative order of  $S_{j+1}$  and  $S_{i+1}$  can be determined from their position in  $SA^{12}$ .
- Case 2: If  $i \bmod 3 = 2$ , we compare the triples  $(s[i], s[i + 1], S_{i+2})$  and  $(s[j], s[j + 1], S_{j+2})$ .

# Recap



# C codes

---

- 50 lines of C++ codes were given in J.C.M. Baeten et al. (Eds.): ICALP 2003, LNCS 2719, pp. 943–955, 2003.
- <http://www.mpi-inf.mpg.de/~sanders/programs/suffix/>