

# Informe 1 2025

## Detección de Heavy Hitters usando Sketches

Ignacio Barría Concha - Nicolás Jarpa Jeldres - Nicolás Pino Leal

Los archivos asociados a este trabajo con el objetivo de identificar kmers canonicos como heavy hitters usando sketches estan en el siguiente enlace:

<https://github.com/Pinox084/Tarea1TMGVD>

### Actividad 1: Preparación de datos y ground truth

#### Objetivo

El objetivo de esta actividad fue definir el conjunto de *k-mers*, calcular un *ground truth* de frecuencias exactas en un subconjunto de datos representativo y establecer umbrales de *heavy hitters (HH)*.

#### Metodología

##### 1. Procesamiento de genomas

- Se trabajó con archivos de genomas en formato **FASTA**.
- Se limpiaron las secuencias dejando únicamente bases **A, C, G, T**, descartando aquellas posiciones con caracteres ambiguos como **N**.
- Cada archivo listado en listaarchivos.txt fue procesado de manera independiente.
- Conservábamos un 1MB de memoria para HH de cada archivo FASTA y los juntamos en un solo archivo.

##### 2. Extracción de k-mers canónicos

- Se implemento un archivo principal llamado kmersheader.hpp para almacenar funciones que permiten codificar los kmers en base 2 y obtener sus formas canónicas.
- La función que entrega los kmers canónicos recibe una sola una línea string de tamaño k que compara el kmer y su complemento y retornar su forma canónica.

- Se consideró tanto el *k-mer* como su reverso complementario, guardando siempre el **canónico** (el menor en orden lexicográfico).

### 3. Cálculo de frecuencias exactas

- Se desarrolló el código que realiza un **conteo exacto de k-mers** utilizando un `unordered_map<u64,uint64_t>`. Donde estas, se obtienen sus frecuencias exactas, del archivo `GTExacto.cpp`.
- El resultado es un archivo con todas las frecuencias exactas (`out_exact.txt`), donde cada línea contiene:

KMER	Frecuencia
...	...
112230021210212110112	10
...	...

Guardados en formato de 2 bits

### 4. Definición de umbral $\phi$ y obtención de heavy hitters

- El umbral se definió como:

$$f(x) \geq \phi N$$

donde  $f(x)$  es la frecuencia de un *k-mer* y  $N$  el total de *k-mers* válidos procesados en el archivo.

- Se probaron valores de  $\phi$ , y coincidimos en el  $\phi$  fijo de  $2e-6$ , esto gracias a ensayo y error en los algoritmos de `countSketch` y `towerSketch` con `CountMin` y `Conservative Update (TS CMCU)`. Esto debido a que en los dos scripts de la actividad 2 el  $\phi$  afectaba mucho más de lo que afectaba en la actividad 1, por tanto, llegamos a una conclusión de que teniendo un  $\phi=2e-6$  dejaría buenos resultados para todo lo pedido.
- El conjunto de heavy hitters de cada archivo se guardó en `HHexactos.txt` y `HHexactos21.txt`, respectivamente para un  $K=31$  y un  $K=21$ .
- Para cumplir con la restricción de memoria, se implementó un límite de 1 MB por archivo, y se junta en los txt resultantes `HHexactos.txt` y `HHexactos21.txt`.

## Resultados

### 1. Frecuencias exactas (Ground Truth)

- Se obtuvo para cada archivo el conteo exacto de todos los *k-mers*.
- Esto constituye el **ground truth** que servirá para evaluar la precisión de los algoritmos de sketch.

### 2. Heavy Hitters

- Para cada archivo se generó la lista de HH exactos con su frecuencia.
- Los resultados se almacenaron en un único archivo de salida HHexactos.txt (para  $k=31$ ) y HHexactos21.txt (para  $k=21$ ), con secciones separadas por archivo.
- Cada bloque contiene como máximo 1 MB de información para cumplir con las restricciones de la actividad.

### 3. Medición de memoria

- Se registró la memoria aproximada utilizada por el `unordered_map` durante el conteo exacto.

## Actividad 2

Se implementaron las estructuras CS y TS-CMCU con funciones `insert(kmer)` y `estimate(kmer)` considerando múltiples filas por tamaño de contadores, a continuación, se presenta el pseudocódigo de ambas implementaciones.

### Pseudocódigo – Count Sketch

Estructura `CountSketch(d, w)`:

`tabla[d][w]  $\leftarrow$  0`

`hashSeeds[j]  $\leftarrow$  valores distintos para cada fila  $j$`

`signSeeds[j]  $\leftarrow$  valores distintos para cada fila  $j$`

Función `SIGN(key, j)`:

`h  $\leftarrow$  Hash(key, signSeeds[j])`

Si ( $h$  es impar) retornar  $+1$

Si ( $h$  es par) retornar  $-1$

Procedimiento UPDATE(key, delta):

Para cada fila  $j = 0 \dots d-1$ :

$h \leftarrow \text{Hash}(\text{key}, \text{hashSeeds}[j])$

$\text{col} \leftarrow h \bmod w$

$s \leftarrow \text{SIGN}(\text{key}, j)$

$\text{tabla}[j][\text{col}] \leftarrow \text{tabla}[j][\text{col}] + \text{delta} * s$

Función ESTIMATE(key):

Para cada fila  $j = 0 \dots d-1$ :

$h \leftarrow \text{Hash}(\text{key}, \text{hashSeeds}[j])$

$\text{col} \leftarrow h \bmod w$

$s \leftarrow \text{SIGN}(\text{key}, j)$

$\text{est}[j] \leftarrow \text{tabla}[j][\text{col}] * s$

retornar mediana( $\text{est}[0..d-1]$ )

update(key, delta)  $\rightarrow$  Inserción / actualización

Para cada una de las “d” filas de la tabla:

1. Se calcula una posición hash col para la clave.
2. Se obtiene un signo aleatorio determinístico +1 o -1 usando otra función hash.
3. Se actualiza el contador de la celda correspondiente con  $\text{delta} * \text{signo}$ .

Esto permite compensar colisiones: si dos claves distintas caen en la misma celda, sus contribuciones pueden tener signos distintos, reduciendo el sesgo.

estimate(key)  $\rightarrow$  Estimación de frecuencia

Para cada fila:

1. Se calcula la misma celda y el mismo signo que se usó en la inserción.
2. Se recupera el valor almacenado y se multiplica por el signo  $\rightarrow$  esto “deshace” el efecto del signo aplicado en update.
3. Se toma la mediana de las  $d$  estimaciones  $\rightarrow$  esto reduce el efecto de posibles colisiones grandes en alguna fila.

El resultado es una estimación no sesgada de la frecuencia de la clave, con alta probabilidad de estar cerca del valor real si “ $d$ ” y “ $w$ ” son lo suficientemente grandes.

## Pseudocódigo -TS-CMCU

Estructura CountMinCU(width, rows):

$\text{tabla}[\text{rows}][\text{width}] \leftarrow 0$

$\text{seeds}[r] \leftarrow$  valores distintos para cada fila  $r$

Procedimiento INSERT(key):

Para cada fila  $r = 0 \dots \text{rows}-1$ :

$h \leftarrow \text{Hash}(\text{key}, \text{seeds}[r])$

$\text{idx}[r] \leftarrow h \bmod \text{width}$

$\text{minv} \leftarrow \infty$

Para cada fila  $r$ :

$\text{minv} \leftarrow \min(\text{minv}, \text{tabla}[r][\text{idx}[r]])$

Para cada fila  $r$ :

Si  $\text{tabla}[r][\text{idx}[r]] == \text{minv}$ :

$\text{tabla}[r][\text{idx}[r]] \leftarrow \text{tabla}[r][\text{idx}[r]] + 1$

Función ESTIMATE(key):

$\text{ans} \leftarrow \infty$

Para cada fila  $r$ :

$h \leftarrow \text{Hash}(\text{key}, \text{seeds}[r])$

```

    idx ← h mod width

    ans ← min(ans, tabla[r][idx])

    Si ans == ∞ retornar 0

    retornar ans

```

insert(key) → Inserción con Conservative Update

1. Se calcula en qué posición (idx) cae la clave en cada una de las rows filas.
2. Se busca el **mínimo valor actual** entre esas posiciones.
3. Se incrementan **solo las celdas que tienen el valor mínimo**.

Esto evita **sobreestimaciones innecesarias**: en el Count-Min clásico se incrementan todas las celdas, lo que puede amplificar errores por colisiones. Con Conservative Update solo se incrementan las celdas más “rezagadas”, manteniendo las otras igual.

estimate(key) → Estimación

1. Se calcula la celda correspondiente en cada fila.
2. Se devuelve el **mínimo de todos los contadores** asociados a la clave.

El mínimo actúa como una cota superior del conteo real, ya que por colisiones el valor en cada celda solo puede **sobreestimar**, nunca subestimar. Tomar el mínimo entre varias filas reduce el error.

Estructura TowerSketch(widths[], rows\_per\_level[], seed0):

Para cada nivel i:

Crear CountMinCU con width = widths[i], rows = rows\_per\_level[i], seed = seed0 + i

Añadirlo a levels

Procedimiento INSERT(key):

Para cada nivel en levels:

nivel.INSERT(key)

Función ESTIMATE(key):

best  $\leftarrow \infty$

Para cada nivel en levels:

e  $\leftarrow$  nivel.ESTIMATE(key)

best  $\leftarrow \min(\text{best}, e)$

Si best ==  $\infty$  retornar 0

retornar best

Tower Sketch combina varios Count-Min Sketch con distintos anchos y números de filas (niveles).

- Cada inserción se propaga por todos los niveles.
- Para estimar, se toma el mínimo entre las estimaciones de todos los niveles.

Esto mejora la precisión y estabilidad respecto a usar un único sketch, especialmente en rangos amplios de frecuencias (muy útiles para heavy hitters o k-mers con distribuciones sesgadas).

## Calibraciones

### Count Sketch:

Se calibra probando distintas combinaciones de número de filas ( $d$ ) y ancho de tabla ( $w$ ). Para cada configuración, se insertan todos los k-mers de un grupo seleccionado de archivos y se comparan las estimaciones con las frecuencias reales, midiendo error (MAE/MRE) y calidad en la detección de heavy hitters (Precisión, Recall, F1). Los parametros utilizados son los siguientes:

vector<int> ds = {9, 10, 11, 15, 20};

vector<int> ws = {10000, 20000, 50000, 100000, 200000, 500000, 1000000, 2000000};

A continuación, se muestra “d” y “w” de prueba para un  $\varphi$  de 2e-06 de tamaños del k-mer de 21.

N°	d	w	Tamaño (bytes)	MAE	MRE	Precisión	Recall	F1
1	9	1 000 000	72 000 000	0.254 749	0.2541 63	0.080717 5	1.000	0.1493 78
2	15	500 000	60 000 000	0.350 795	0.3500 06	0.038095 2	0.667	0.0720 721
3	20	2 000 000	320 000 000	0.022 5919	0.0225 223	0.18	1.000	0.3050 85
4	10	100 000	8 000 000	1.248 64	1.2458 6	8.51396e -05	0.667	0.0001 703
5	11	500 000	44 000 000	0.432 045	0.4310 76	0.024163 6	0.722	0.0467 626
6	9	200 000	14 400 000	0.858 456	0.8565 42	0.000476 304	0.556	0.0009 518
7	20	1 000 000	160 000 000	0.107 176	0.1069 25	0.130435	1.000	0.2307 69
8	15	1 000 000	120 000 000	0.147 619	0.1472 80	0.112583	0.944	0.2011 83
9	10	500 000	40 000 000	0.490 636	0.4895 26	0.009451 8	0.833	0.0186 916
10	9	50 000	3 600 000	1.786 61	1.7826 7	2.92946e -05	0.333	5.8584 e-05



A continuación, se muestra “d” y “w” de prueba para un  $\varphi$  de 2e-06 de tamaños del k-mer de 31.

N°	d	w	Tamaño (bytes)	MAE	MRE	Precis ión	Recall	F1
1	20	500 000	80 000 000	0.2401 39	0.239809	0.003 48172	0.875	0.00693 584
2	9	500 000	36 000 000	0.4237 60	0.423208	0.000 96125 4	0.8125	0.00192 024
3	11	200 000	17 600 000	0.6946 83	0.693736	0.000 20577 6	0.875	0.00041 146
4	15	20 000	2 400 000	2.0172 30	2.014600	0.000 02558 6	0.5625	0.00005 117
5	9	1 000 000	72 000 000	0.1995 86	0.199306	0.005 50863 0	0.9375	0.01095 290
6	11	10 000	880 000	3.3141 40	3.309930	0.000 01988 9	0.5625	0.00003 978
7	20	1 000 000	160 000 000	0.0718 52	0.071729	0.011 51080 0	1.0000	0.02275 960
8	15	500 000	60 000 000	0.2944 62	0.294074	0.003 39751 0	0.9375	0.00677 048
9	10	10 000	800 000	3.6521 50	3.647460	0.000 01627 4	0.625	0.00003 255
10	11	1 000 000	88 000 000	0.1587 30	0.158496	0.007 81250 0	1.0000	0.01550 390

### Tower Sketch:

Se calibra definiendo varios niveles con anchos decrecientes y filas fijas, ajustando los parámetros según  $\phi$ ,  $\delta$  y la memoria disponible. Se evalúa igual que Count Sketch, tomando métricas de error y heavy hitters para comparar configuraciones. Los parámetros utilizados fueron los siguientes

```
vector<vector<size_t>> configs_widths = {  
  
    {1<<20, 1<<19, 1<<18, 1<<17, 1<<16},  
  
    {1<<21, 1<<20, 1<<19, 1<<18, 1<<17},  
  
    {1<<23, 1<<22, 1<<21, 1<<20, 1<<19},  
  
    {1<<24, 1<<23, 1<<22, 1<<21, 1<<20, 1<<19, 1<<18}  
}  
  
vector<vector<size_t>> configs_rows = {  
  
    {5, 5, 4, 4, 3},  
  
    {6, 5, 5, 4, 4},  
  
    {7, 6, 6, 5, 5},  
  
    {10, 8, 7, 7, 6, 6, 5}  
}
```

A continuación, se muestra  $d$  y  $w$  de prueba para un  $\varphi$  de  $2e-06$  de tamaños del k-mer de 21.

Co nfi g	N	MAE	RMS E	Max Abs Error	Max Rel Error	Precision	Recall	F1Score
cfg 1	1302461	0.004 09202	0.074 2833	12	12	0.323442	1	0.488789
cfg 3	1514500	0.013 299	0.121 375	6	6	0.188119	1	0.316667
cfg 2	3982388	0.018 5966	0.209 265	35	35	0.390465	1	0.561633

cfg 0	2277733	0.174 591	0.454 005	29	29	0.460581	1	0.630682
cfg 1	1942545	0.009 11808	0.251 977	41	41	0.240964	1	0.38835
cfg 2	1580689	0.002 05053	0.047 1459	7	7	0.638889	1	0.779661
cfg 1	3980036	0.100 015	0.372 99	76	76	0.0797721	1	0.147757
cfg 2	2284665	0.006 94263	0.098 1644	12	12	0.0740741	1	0.137931
cfg 0	1514500	0.053 1368	0.234 341	6	6	0.180952	1	0.306452
cfg 3	3982388	0.018 5958	0.209 263	35	35	0.390465	1	0.561633

A continuación, se muestra d y w de prueba para un  $\varphi$  de 2e-06 de tamaños del k-mer de 31.

<b>Co nfi g</b>	<b>N</b>	<b>MAE</b>	<b>RMSE</b>	<b>Max Abs Error</b>	<b>Max Rel Error</b>	<b>Precision</b>	<b>Reca ll</b>	<b>F1Score</b>
cfg 2	1260 958	0.012 9983	0.1198 78	7	7	0.0264026	1	0.0514469
cfg 1	1084 253	0.003 7918	0.0759 957	9	9	0.186047	1	0.313725
cfg 3	1316 363	0.002 0622 7	0.0472 008	5	5	0.121212	1	0.216216
cfg 0	3318 223	0.454 088	0.7172 79	29	29	0.193969	1	0.324915
cfg 1	1897 122	0.018 3609	0.2808 1	25	25	0.209479	1	0.346395

cfg 3	1903 784	0.006 1672 5	0.0950 366	11	11	0.278607	1	0.435798
cfg 0	1618 601	0.063 3776	0.4062 76	35	35	0.0786651	1	0.145856
cfg 0	7978 64	0.002 8542 6	0.0545 628	2	2	0.0462527	1	0.0884159
cfg 1	3316 461	0.058 5379	0.3212 62	52	52	0.0018832 4	1	0.0037594
cfg 2	3318 223	0.018 6584	0.2329 47	29	29	0.199664	1	0.332867

### Actividad 3: Detección de Heavy Hitters y Validación

Se ejecutan las dos estructuras usando los parametros con mejor rendimiento que se obtuvieron tras la calibración, para la detección de heavy hitters se usaron los siguientes parametros:

- $\Phi = 2e-6$  (0.000002)
- Count Sketch:  $d = 11$ ,  $width = 2.000.000$
- TS-CMCU:  $widths = \{1 < 22, 1 < 21, 1 < 20, 1 < 19, 1 < 18\}$ ,  $rows = \{6, 6, 5, 5, 4\}$

TS-CMCU con  $k = 21$

id x	name	TP	FP	F N	precis ion	rec all	f1	hh_ex act	hh_ est
1	GCA_006152045.1_ASM615204v1_genomic.fna	1	52 4	0	0,001 905	1	0,003 802	1	525
2	GCA_018421455.1_ASM1842145v1_genomic.fna	23 8	95 4	0	0,199 664	1	0,332 867	238	119 2
3	GCA_020118255.1_ASM2011825v1_genomic.fna	56	11 72	0	0,045 603	1	0,087 227	56	122 8
4	GCA_021919605.1_PDT001092240.1_genomic.fna	42	91 6	0	0,043 841	1	0,084	42	958
5	GCA_021953145.1_PDT001020446.1_genomic.fna	17	74 6	0	0,022 28	1	0,043 59	17	763
6	GCA_021972575.1_PDT001013406.1_genomic.fna	14	79 5	0	0,017 305	1	0,034 022	14	809
7	GCA_022062785.1_PDT000876120.1_genomic.fna	30	77 6	0	0,037 221	1	0,071 77	30	806

8	GCA_023315275.2_PDT001299634 .2_genomic.fna	17	90 4	0	0,018 458	1	0,036 247	17	921
9	GCA_024452925.1_PDT001370221 .1_genomic.fna	55	92 7	0	0,056 008	1	0,106 075	55	982
10	GCA_024732165.1_PDT001378927 .1_genomic.fna	44	82 5	0	0,050 633	1	0,096 386	44	869
11	GCA_026006075.1_ASM2600607v 1_genomic.fna	49 5	10 04	0	0,330 22	1	0,496 489	495	149 9
12	GCA_026006095.1_ASM2600609v 1_genomic.fna	53 3	95 4	0	0,358 44	1	0,527 723	533	148 7
13	GCA_026305215.1_PDT001493239 .1_genomic.fna	1	55 5	0	0,001 799	1	0,003 591	1	556
14	GCA_031045075.2_PDT001744190 .2_genomic.fna	53	91 0	0	0,055 036	1	0,104 331	53	963
15	GCA_032567175.1_ASM3256717v 1_genomic.fna	48	85 5	0	0,053 156	1	0,100 946	48	903
16	GCA_033106465.1_PDT001960943 .1_genomic.fna	56	94 0	0	0,056 225	1	0,106 464	56	996
17	GCA_037052065.1_ASM3705206v 1_genomic.fna	0	41 1	0	0	0	0	0	411
18	GCA_037203645.1_ASM3720364v 1_genomic.fna	11 1	94 3	0	0,105 313	1	0,190 558	111	105 4
19	GCA_043678295.1_ASM4367829v 1_genomic.fna	80	93 7	0	0,078 663	1	0,145 852	80	101 7
20	GCA_043678335.1_ASM4367833v 1_genomic.fna	30 6	98 0	0	0,237 947	1	0,384 422	306	128 6
21	GCA_043678355.1_ASM4367835v 1_genomic.fna	60 3	89 5	0	0,402 537	1	0,574 012	603	149 8
22	GCA_043678375.1_ASM4367837v 1_genomic.fna	48 3	10 00	0	0,325 691	1	0,491 353	483	148 3
23	GCA_043950085.1_ASM4395008v 1_genomic.fna	20	55 5	0	0,034 783	1	0,067 227	20	575
24	GCA_943323015.2_HRS-ES2-bin- 440_genomic.fna	57 6	34 54	0	0,142 928	1	0,250 109	576	403 0
25	GCA_943323415.2_HRS-ES4-bin- 10_genomic.fna	23 66	36 22	0	0,395 124	1	0,566 435	2366	598 8
26	GCA_943323865.2_HRS-ES8-bin- 118_genomic.fna	50 4	33 80	0	0,129 763	1	0,229 717	504	388 4
27	GCA_943325355.2_HRSG-E10-bin- 12_genomic.fna	8	65 6	0	0,012 048	1	0,023 81	8	664
28	GCA_943326675.2_HRSG-E12-bin- 108_genomic.fna	71	66 4	0	0,096 599	1	0,176 179	71	735
29	GCA_943327095.2_HRSG-E11-bin- 71_genomic.fna	12	40 3	0	0,028 916	1	0,056 206	12	415
30	GCA_943327695.2_HRS-HS3-bin- 186_genomic.fna	1	10 6	0	0,009 346	1	0,018 519	1	107

3 1	GCA_943329525.2_HRS-HS6-bin-101_genomic.fna	10 8	11 31	0	0,087 167	1	0,160 356	108	123 9
3 2	GCA_943329805.2_HRSG-H13-bin-193_genomic.fna	44 2	16 68	0	0,209 479	1	0,346 395	442	211 0
3 3	GCA_943331045.1_HRSG-E2-bin-118_genomic.fna	16	11 6	0	0,121 212	1	0,216 216	16	132
3 4	GCA_943331235.2_HRSG-H12-bin-94_genomic.fna	24	49 6	0	0,046 154	1	0,088 235	24	520
3 5	GCA_943332185.1_HRS-HS9-bin-103_genomic.fna	67	12 10	0	0,052 467	1	0,099 702	67	127 7
3 6	GCA_943332625.1_HRS-HS9-bin-149_genomic.fna	11	15 7	0	0,065 476	1	0,122 905	11	168
3 7	GCA_943333245.1_HRSG-E1-bin-8_genomic.fna	16 73	85 81	0	0,163 156	1	0,280 54	1673	102 54
3 8	GCA_943333305.2_HRSG-H6-bin-46_genomic.fna	26	26 13	0	0,009 852	1	0,019 512	26	263 9
3 9	GCA_943334125.2_HRSG-E6-bin-23_genomic.fna	14 4	16 02	0	0,082 474	1	0,152 381	144	174 6
4 0	GCA_943334475.2_HRSG-H3-bin-66_genomic.fna	18 8	46 02	0	0,039 248	1	0,075 532	188	479 0
4 1	GCA_943334715.2_HRSG-H6-bin-64_genomic.fna	56	14 5	0	0,278 607	1	0,435 798	56	201
4 2	GCA_943334795.2_HRSG-H6-bin-28_genomic.fna	48	21 0	0	0,186 047	1	0,313 725	48	258
4 3	GCA_945860935.1_GE-03apr19-182_genomic.fna	0	34 4	0	0	0	0	0	344
4 4	GCA_945867275.1_TH-11nov19-195_genomic.fna	41	69 9	0	0,055 405	1	0,104 994	41	740
4 5	GCA_945870865.1_TrH-03may19-115_genomic.fna	10 8	10 14	0	0,096 257	1	0,175 61	108	112 2
4 6	GCA_945870955.1_MoE-23oct19-325_genomic.fna	66	75 4	0	0,080 488	1	0,148 984	66	820
4 7	GCA_945872895.1_MaE-04nov19-161_genomic.fna	16	59 0	0	0,026 403	1	0,051 447	16	606
4 8	GCA_945874295.1_AH-24oct19-10_genomic.fna	25	15 0	0	0,142 857	1	0,25	25	175
4 9	GCA_945877175.1_ZE-13nov19-95_genomic.fna	17 4	34 10	0	0,048 549	1	0,092 602	174	358 4
5 0	GCA_945902215.1_MoH-23oct19-118_genomic.fna	20 5	14 72	0	0,122 242	1	0,217 853	205	167 7

Count Sketch con  $k = 21$

id x	name	TP	FP	F N	precis ion	recall	f1	hh_e xact	hh_ est
1	GCA_006152045.1_ASM615204v1_genomic.fna	1	515	0	0,001938	1	0,003868	1	516
2	GCA_018421455.1_ASM1842145v1_genomic.fna	231	973	7	0,19186	0,970588	0,320388	238	1204
3	GCA_020118255.1_ASM2011825v1_genomic.fna	53	1190	3	0,042639	0,946429	0,081601	56	1243
4	GCA_021919605.1_PDT001092240.1_genomic.fna	42	925	0	0,043433	1	0,083251	42	967
5	GCA_021953145.1_PDT001020446.1_genomic.fna	17	726	0	0,02288	1	0,044737	17	743
6	GCA_021972575.1_PDT001013406.1_genomic.fna	13	762	1	0,016774	0,928571	0,032953	14	775
7	GCA_022062785.1_PDT000876120.1_genomic.fna	29	780	1	0,035847	0,966667	0,06913	30	809
8	GCA_023315275.2_PDT001299634.2_genomic.fna	17	924	0	0,018066	1	0,035491	17	941
9	GCA_024452925.1_PDT001370221.1_genomic.fna	54	898	1	0,056723	0,981818	0,107249	55	952
10	GCA_024732165.1_PDT001378927.1_genomic.fna	40	818	4	0,04662	0,909091	0,088692	44	858
11	GCA_026006075.1_ASM2600607v1_genomic.fna	463	1021	32	0,311995	0,935354	0,467913	495	1484
12	GCA_026006095.1_ASM2600609v1_genomic.fna	500	992	33	0,335121	0,938086	0,493827	533	1492
13	GCA_026305215.1_PDT001493239.1_genomic.fna	1	541	0	0,001845	1	0,003683	1	542
14	GCA_031045075.2_PDT001744190.2_genomic.fna	50	905	3	0,052356	0,943396	0,099206	53	955
15	GCA_032567175.1_ASM3256717v1_genomic.fna	46	880	2	0,049676	0,958333	0,094456	48	926
16	GCA_033106465.1_PDT001960943.1_genomic.fna	51	915	5	0,052795	0,910714	0,099804	56	966
17	GCA_037052065.1_ASM3705206v1_genomic.fna	0	428	0	0	0	0	0	428
18	GCA_037203645.1_ASM3720364v1_genomic.fna	102	953	9	0,096682	0,918919	0,174957	111	1055
19	GCA_043678295.1_ASM4367829v1_genomic.fna	77	1048	3	0,068444	0,9625	0,127801	80	1125
20	GCA_043678335.1_ASM4367833v1_genomic.fna	274	1115	32	0,197264	0,895425	0,323304	306	1389
21	GCA_043678355.1_ASM4367835v1_genomic.fna	557	1020	46	0,353202	0,923715	0,511009	603	1577

2	GCA_043678375.1_ASM4367837	44	114	4	0,279	0,915	0,428		158
2	v1_genomic.fna	2	0	1	393	114	087	483	2
2	GCA_043950085.1_ASM4395008				0,031		0,060		
3	v1_genomic.fna	18	556	2	359	0,9	606	20	574
2	GCA_943323015.2_HRS-ES2-bin-	57	631		0,083	0,994	0,153		689
4	440_genomic.fna	3	8	3	152	792	475	576	1
2	GCA_943323415.2_HRS-ES4-bin-	23	666		0,261	0,996	0,414		902
5	10_genomic.fna	58	4	8	361	619	12	2366	2
2	GCA_943323865.2_HRS-ES8-bin-	50	689		0,067	0,998	0,127		740
6	118_genomic.fna	3	7	1	973	016	277	504	0
2	GCA_943325355.2_HRSG-E10-				0,009		0,019		
7	bin-12_genomic.fna	8	826	0	592	1	002	8	834
2	GCA_943326675.2_HRSG-E12-				0,080	0,957	0,147		
8	bin-108_genomic.fna	68	781	3	094	746	826	71	849
2	GCA_943327095.2_HRSG-E11-				0,014	0,916	0,028		
9	bin-71_genomic.fna	11	762	1	23	667	025	12	773
3	GCA_943327695.2_HRS-HS3-bin-				0,008		0,017		
0	186_genomic.fna	1	111	0	929	1	699	1	112
3	GCA_943329525.2_HRS-HS6-bin-	10	334		0,031	0,990	0,060		345
1	101_genomic.fna	7	4	1	006	741	129	108	1
3	GCA_943329805.2_HRSG-H13-	44	170		0,205	0,995	0,340		214
2	bin-193_genomic.fna	0	2	2	415	475	557	442	2
3	GCA_943331045.1_HRSG-E2-bin-				0,099		0,180		
3	118_genomic.fna	16	145	0	379	1	791	16	161
3	GCA_943331235.2_HRSG-H12-				0,042		0,081		
4	bin-94_genomic.fna	24	539	0	629	1	772	24	563
3	GCA_943332185.1_HRS-HS9-bin-		135		0,046	0,985	0,088		142
5	103_genomic.fna	66	6	1	414	075	65	67	2
3	GCA_943332625.1_HRS-HS9-bin-				0,062		0,117		
6	149_genomic.fna	11	165	0	5	1	647	11	176
3	GCA_943333245.1_HRSG-E1-bin-	16	102		0,140	0,998	0,246		118
7	8_genomic.fna	71	12	2	621	805	533	1673	83
3	GCA_943333305.2_HRSG-H6-bin-		262		0,009	0,961	0,018		265
8	46_genomic.fna	25	7	1	427	538	671	26	2
3	GCA_943334125.2_HRSG-E6-bin-	14	163		0,080	0,993	0,148		177
9	23_genomic.fna	3	6	1	382	056	726	144	9
4	GCA_943334475.2_HRSG-H3-bin-	18	466		0,038		0,074		485
0	66_genomic.fna	8	8	0	715	1	544	188	6
4	GCA_943334715.2_HRSG-H6-bin-				0,252	0,982	0,401		
1	64_genomic.fna	55	163	1	294	143	46	56	218
4	GCA_943334795.2_HRSG-H6-bin-				0,177		0,300		
2	28_genomic.fna	48	223	0	122	1	94	48	271
4	GCA_945860935.1_GE-03apr19-								
3	182_genomic.fna	0	379	0	0	0	0	0	379
4	GCA_945867275.1_TH-11nov19-				0,053		0,100		
4	195_genomic.fna	41	730	0	178	1	985	41	771



4 5	GCA_945870865.1_TrH- 03may19-115_genomic.fna	10 8	110 9	0	0,088 743	1	0,163 019	108	121 7
4 6	GCA_945870955.1_MoE- 23oct19-325_genomic.fna	66	768	0	0,079 137	1	0,146 667	66	834
4 7	GCA_945872895.1_MaE- 04nov19-161_genomic.fna	16	837	0	0,018 757	1	0,036 824	16	853
4 8	GCA_945874295.1_AH-24oct19- 10_genomic.fna	25	157	0	0,137 363	1	0,241 546	25	182
4 9	GCA_945877175.1_ZE-13nov19- 95_genomic.fna	17 0	343 2	4	0,047 196	0,977 011	0,090 042	174	360 2
5 0	GCA_945902215.1_MoH- 23oct19-118_genomic.fna	20 3	149 6	2	0,119 482	0,990 244	0,213 235	205	169 9

TS-CMCU con k= 31

id x	name	TP	FP	F N	precis ion	rec all	f1	hh_ex act	hh_ est
1	GCA_006152045.1_ASM615204v1 _genomic.fna	1	52 4	0	0,001 905	1	0,003 802	1	525
2	GCA_018421455.1_ASM1842145v 1_genomic.fna	23 8	95 4	0	0,199 664	1	0,332 867	238	119 2
3	GCA_020118255.1_ASM2011825v 1_genomic.fna	56	11 72	0	0,045 603	1	0,087 227	56	122 8
4	GCA_021919605.1_PDT001092240 .1_genomic.fna	42	91 6	0	0,043 841	1	0,084	42	958
5	GCA_021953145.1_PDT001020446 .1_genomic.fna	17	74 6	0	0,022 28	1	0,043 59	17	763
6	GCA_021972575.1_PDT001013406 .1_genomic.fna	14	79 5	0	0,017 305	1	0,034 022	14	809
7	GCA_022062785.1_PDT000876120 .1_genomic.fna	30	77 6	0	0,037 221	1	0,071 77	30	806
8	GCA_023315275.2_PDT001299634 .2_genomic.fna	17	90 4	0	0,018 458	1	0,036 247	17	921
9	GCA_024452925.1_PDT001370221 .1_genomic.fna	55	92 7	0	0,056 008	1	0,106 075	55	982
1 0	GCA_024732165.1_PDT001378927 .1_genomic.fna	44	82 5	0	0,050 633	1	0,096 386	44	869
1 1	GCA_026006075.1_ASM2600607v 1_genomic.fna	49 5	10 04	0	0,330 22	1	0,496 489	495	149 9
1 2	GCA_026006095.1_ASM2600609v 1_genomic.fna	53 3	95 4	0	0,358 44	1	0,527 723	533	148 7
1 3	GCA_026305215.1_PDT001493239 .1_genomic.fna	1	55 5	0	0,001 799	1	0,003 591	1	556
1 4	GCA_031045075.2_PDT001744190 .2_genomic.fna	53	91 0	0	0,055 036	1	0,104 331	53	963

1 5	GCA_032567175.1_ASM3256717v 1_genomic.fna	48	85 5	0	0,053 156	1	0,100 946	48	903
1 6	GCA_033106465.1_PDT001960943 .1_genomic.fna	56	94 0	0	0,056 225	1	0,106 464	56	996
1 7	GCA_037052065.1_ASM3705206v 1_genomic.fna	0	41 1	0	0	0	0	0	411
1 8	GCA_037203645.1_ASM3720364v 1_genomic.fna	11 1	94 3	0	0,105 313	1	0,190 558	111	105 4
1 9	GCA_043678295.1_ASM4367829v 1_genomic.fna	80	93 7	0	0,078 663	1	0,145 852	80	101 7
2 0	GCA_043678335.1_ASM4367833v 1_genomic.fna	30 6	98 0	0	0,237 947	1	0,384 422	306	128 6
2 1	GCA_043678355.1_ASM4367835v 1_genomic.fna	60 3	89 5	0	0,402 537	1	0,574 012	603	149 8
2 2	GCA_043678375.1_ASM4367837v 1_genomic.fna	48 3	10 00	0	0,325 691	1	0,491 353	483	148 3
2 3	GCA_043950085.1_ASM4395008v 1_genomic.fna	20	55 5	0	0,034 783	1	0,067 227	20	575
2 4	GCA_943323015.2_HRS-ES2-bin- 440_genomic.fna	57 6	34 54	0	0,142 928	1	0,250 109	576	403 0
2 5	GCA_943323415.2_HRS-ES4-bin- 10_genomic.fna	23 66	36 22	0	0,395 124	1	0,566 435	2366	598 8
2 6	GCA_943323865.2_HRS-ES8-bin- 118_genomic.fna	50 4	33 80	0	0,129 763	1	0,229 717	504	388 4
2 7	GCA_943325355.2_HRSG-E10-bin- 12_genomic.fna	8	65 6	0	0,012 048	1	0,023 81	8	664
2 8	GCA_943326675.2_HRSG-E12-bin- 108_genomic.fna	71	66 4	0	0,096 599	1	0,176 179	71	735
2 9	GCA_943327095.2_HRSG-E11-bin- 71_genomic.fna	12	40 3	0	0,028 916	1	0,056 206	12	415
3 0	GCA_943327695.2_HRS-HS3-bin- 186_genomic.fna	1	10 6	0	0,009 346	1	0,018 519	1	107
3 1	GCA_943329525.2_HRS-HS6-bin- 101_genomic.fna	10 8	11 31	0	0,087 167	1	0,160 356	108	123 9
3 2	GCA_943329805.2_HRSG-H13-bin- 193_genomic.fna	44 2	16 68	0	0,209 479	1	0,346 395	442	211 0
3 3	GCA_943331045.1_HRSG-E2-bin- 118_genomic.fna	16	11 6	0	0,121 212	1	0,216 216	16	132
3 4	GCA_943331235.2_HRSG-H12-bin- 94_genomic.fna	24	49 6	0	0,046 154	1	0,088 235	24	520
3 5	GCA_943332185.1_HRS-HS9-bin- 103_genomic.fna	67	12 10	0	0,052 467	1	0,099 702	67	127 7
3 6	GCA_943332625.1_HRS-HS9-bin- 149_genomic.fna	11	15 7	0	0,065 476	1	0,122 905	11	168
3 7	GCA_943333245.1_HRSG-E1-bin- 8_genomic.fna	16 73	85 81	0	0,163 156	1	0,280 54	1673	102 54

3 8	GCA_943333305.2_HRSG-H6-bin-46_genomic.fna	26	26 13	0	0,009 852	1	0,019 512	26	263 9
3 9	GCA_943334125.2_HRSG-E6-bin-23_genomic.fna	14 4	16 02	0	0,082 474	1	0,152 381	144	174 6
4 0	GCA_943334475.2_HRSG-H3-bin-66_genomic.fna	18 8	46 02	0	0,039 248	1	0,075 532	188	479 0
4 1	GCA_943334715.2_HRSG-H6-bin-64_genomic.fna	56	14 5	0	0,278 607	1	0,435 798	56	201
4 2	GCA_943334795.2_HRSG-H6-bin-28_genomic.fna	48	21 0	0	0,186 047	1	0,313 725	48	258
4 3	GCA_945860935.1_GE-03apr19-182_genomic.fna	0	34 4	0	0	0	0	0	344
4 4	GCA_945867275.1_TH-11nov19-195_genomic.fna	41	69 9	0	0,055 405	1	0,104 994	41	740
4 5	GCA_945870865.1_TrH-03may19-115_genomic.fna	10 8	10 14	0	0,096 257	1	0,175 61	108	112 2
4 6	GCA_945870955.1_MoE-23oct19-325_genomic.fna	66	75 4	0	0,080 488	1	0,148 984	66	820
4 7	GCA_945872895.1_MaE-04nov19-161_genomic.fna	16	59 0	0	0,026 403	1	0,051 447	16	606
4 8	GCA_945874295.1_AH-24oct19-10_genomic.fna	25	15 0	0	0,142 857	1	0,25	25	175
4 9	GCA_945877175.1_ZE-13nov19-95_genomic.fna	17 4	34 10	0	0,048 549	1	0,092 602	174	358 4
5 0	GCA_945902215.1_MoH-23oct19-118_genomic.fna	20 5	14 72	0	0,122 242	1	0,217 853	205	167 7

Count Sketch con  $k = 31$

id x	name	TP	FP	F N	precis ion	recall	f1	hh_e xact	hh_ est
1	GCA_006152045.1_ASM615204v1_genomic.fna	1	515	0	0,001 938	1	0,003 868	1	516
2	GCA_018421455.1_ASM1842145v1_genomic.fna	23 1	973	7	0,191 86	0,970 588	0,320 388	238	120 4
3	GCA_020118255.1_ASM2011825v1_genomic.fna	53	119 0	3	0,042 639	0,946 429	0,081 601	56	124 3
4	GCA_021919605.1_PDT001092240.1_genomic.fna	42	925	0	0,043 433	1	0,083 251	42	967
5	GCA_021953145.1_PDT001020446.1_genomic.fna	17	726	0	0,022 88	1	0,044 737	17	743
6	GCA_021972575.1_PDT001013406.1_genomic.fna	13	762	1	0,016 774	0,928 571	0,032 953	14	775
7	GCA_022062785.1_PDT000876120.1_genomic.fna	29	780	1	0,035 847	0,966 667	0,069 13	30	809

8	GCA_023315275.2_PDT0012996 34.2_genomic.fna	17	924	0	0,018 066	1	0,035 491	17	941
9	GCA_024452925.1_PDT0013702 21.1_genomic.fna	54	898	1	0,056 723	0,981 818	0,107 249	55	952
1 0	GCA_024732165.1_PDT0013789 27.1_genomic.fna	40	818	4	0,046 62	0,909 091	0,088 692	44	858
1 1	GCA_026006075.1_ASM2600607 v1_genomic.fna	46 3	102 1	3 2	0,311 995	0,935 354	0,467 913	495	148 4
1 2	GCA_026006095.1_ASM2600609 v1_genomic.fna	50 0	992	3 3	0,335 121	0,938 086	0,493 827	533	149 2
1 3	GCA_026305215.1_PDT0014932 39.1_genomic.fna	1	541	0	0,001 845	1	0,003 683	1	542
1 4	GCA_031045075.2_PDT0017441 90.2_genomic.fna	50	905	3	0,052 356	0,943 396	0,099 206	53	955
1 5	GCA_032567175.1_ASM3256717 v1_genomic.fna	46	880	2	0,049 676	0,958 333	0,094 456	48	926
1 6	GCA_033106465.1_PDT0019609 43.1_genomic.fna	51	915	5	0,052 795	0,910 714	0,099 804	56	966
1 7	GCA_037052065.1_ASM3705206 v1_genomic.fna	0	428	0	0	0	0	0	428
1 8	GCA_037203645.1_ASM3720364 v1_genomic.fna	10 2	953	9	0,096 682	0,918 919	0,174 957	111	105 5
1 9	GCA_043678295.1_ASM4367829 v1_genomic.fna	77	104 8	3	0,068 444	0,962 5	0,127 801	80	112 5
2 0	GCA_043678335.1_ASM4367833 v1_genomic.fna	27 4	111 5	3 2	0,197 264	0,895 425	0,323 304	306	138 9
2 1	GCA_043678355.1_ASM4367835 v1_genomic.fna	55 7	102 0	4 6	0,353 202	0,923 715	0,511 009	603	157 7
2 2	GCA_043678375.1_ASM4367837 v1_genomic.fna	44 2	114 0	4 1	0,279 393	0,915 114	0,428 087	483	158 2
2 3	GCA_043950085.1_ASM4395008 v1_genomic.fna	18	556	2	0,031 359	0,9	0,060 606	20	574
2 4	GCA_943323015.2_HRS-ES2-bin- 440_genomic.fna	57 3	631 8	3	0,083 152	0,994 792	0,153 475	576	689 1
2 5	GCA_943323415.2_HRS-ES4-bin- 10_genomic.fna	23 58	666 4	8	0,261 361	0,996 619	0,414 12	2366	902 2
2 6	GCA_943323865.2_HRS-ES8-bin- 118_genomic.fna	50 3	689 7	1	0,067 973	0,998 016	0,127 277	504	740 0
2 7	GCA_943325355.2_HRSG-E10- bin-12_genomic.fna	8	826	0	0,009 592	1	0,019 002	8	834
2 8	GCA_943326675.2_HRSG-E12- bin-108_genomic.fna	68	781	3	0,080 094	0,957 746	0,147 826	71	849
2 9	GCA_943327095.2_HRSG-E11- bin-71_genomic.fna	11	762	1	0,014 23	0,916 667	0,028 025	12	773
3 0	GCA_943327695.2_HRS-HS3-bin- 186_genomic.fna	1	111	0	0,008 929	1	0,017 699	1	112

3 1	GCA_943329525.2_HRS-HS6-bin-101_genomic.fna	10 7	334 4	1	0,031 006	0,990 741	0,060 129	108	345 1
3 2	GCA_943329805.2_HRSG-H13-bin-193_genomic.fna	44 0	170 2	2	0,205 415	0,995 475	0,340 557	442	214 2
3 3	GCA_943331045.1_HRSG-E2-bin-118_genomic.fna	16	145	0	0,099 379	1	0,180 791	16	161
3 4	GCA_943331235.2_HRSG-H12-bin-94_genomic.fna	24	539	0	0,042 629	1	0,081 772	24	563
3 5	GCA_943332185.1_HRS-HS9-bin-103_genomic.fna	66	135 6	1	0,046 414	0,985 075	0,088 65	67	142 2
3 6	GCA_943332625.1_HRS-HS9-bin-149_genomic.fna	11	165	0	0,062 5	1	0,117 647	11	176
3 7	GCA_943333245.1_HRSG-E1-bin-8_genomic.fna	16 71	102 12	2	0,140 621	0,998 805	0,246 533	1673	118 83
3 8	GCA_943333305.2_HRSG-H6-bin-46_genomic.fna	25	262 7	1	0,009 427	0,961 538	0,018 671	26	265 2
3 9	GCA_943334125.2_HRSG-E6-bin-23_genomic.fna	14 3	163 6	1	0,080 382	0,993 056	0,148 726	144	177 9
4 0	GCA_943334475.2_HRSG-H3-bin-66_genomic.fna	18 8	466 8	0	0,038 715	1	0,074 544	188	485 6
4 1	GCA_943334715.2_HRSG-H6-bin-64_genomic.fna	55	163	1	0,252 294	0,982 143	0,401 46	56	218
4 2	GCA_943334795.2_HRSG-H6-bin-28_genomic.fna	48	223	0	0,177 122	1	0,300 94	48	271
4 3	GCA_945860935.1_GE-03apr19-182_genomic.fna	0	379	0	0	0	0	0	379
4 4	GCA_945867275.1_TH-11nov19-195_genomic.fna	41	730	0	0,053 178	1	0,100 985	41	771
4 5	GCA_945870865.1_TrH-03may19-115_genomic.fna	10 8	110 9	0	0,088 743	1	0,163 019	108	121 7
4 6	GCA_945870955.1_MoE-23oct19-325_genomic.fna	66	768	0	0,079 137	1	0,146 667	66	834
4 7	GCA_945872895.1_MaE-04nov19-161_genomic.fna	16	837	0	0,018 757	1	0,036 824	16	853
4 8	GCA_945874295.1_AH-24oct19-10_genomic.fna	25	157	0	0,137 363	1	0,241 546	25	182
4 9	GCA_945877175.1_ZE-13nov19-95_genomic.fna	17 0	343 2	4	0,047 196	0,977 011	0,090 042	174	360 2
5 0	GCA_945902215.1_MoH-23oct19-118_genomic.fna	20 3	149 6	2	0,119 482	0,990 244	0,213 235	205	169 9

Gráfico 1: Comparación entre los heavy hitters de k-mers 21 con  $\phi 2e-06$

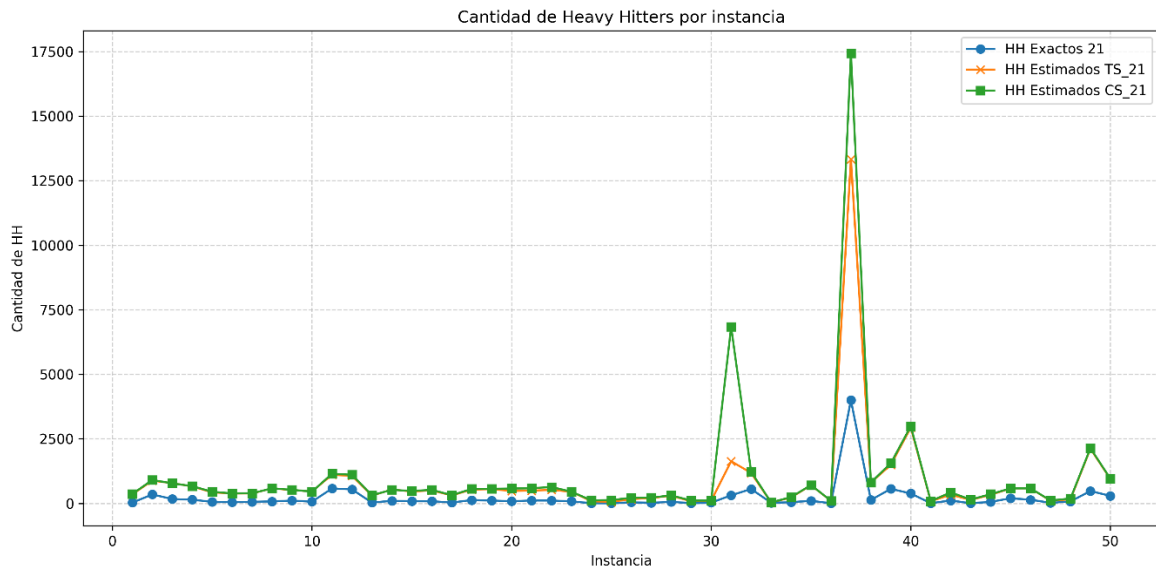
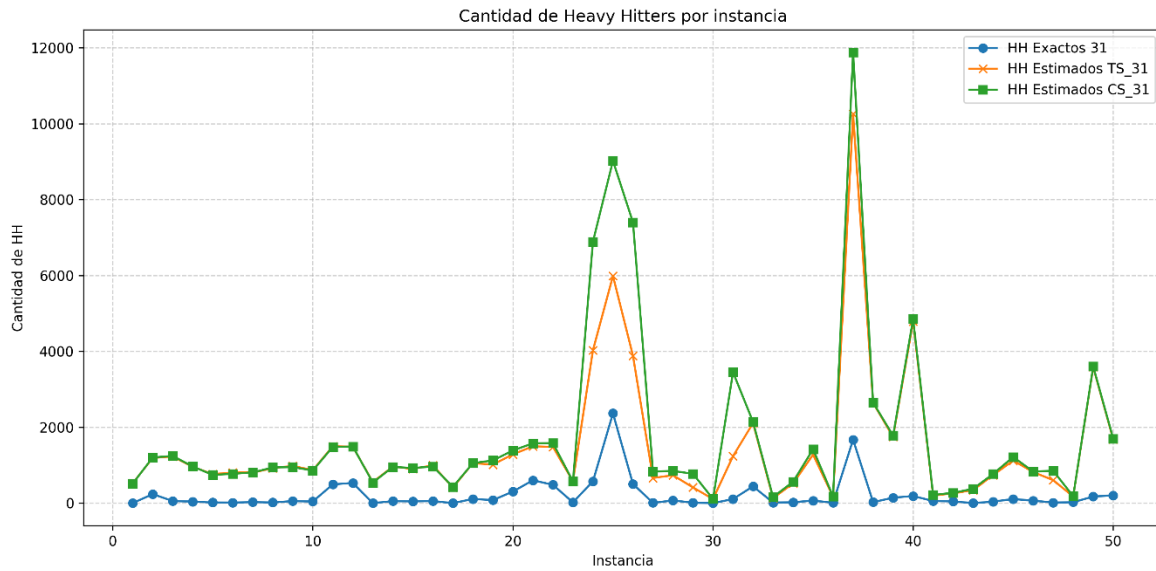


Gráfico 2: Comparación entre los heavy hitters de k-mers 31 con  $\phi 2e-06$



# Análisis comparativo

## Estructura y funcionamiento

- CountSketch (CS):
  - Usa funciones hash y un conjunto de contadores con signo (+1, -1).
  - La estimación se obtiene tomando la mediana de los valores acumulados.
  - Permite compensar ruido y manejar bien distribuciones sesgadas, reduciendo el sesgo de las colisiones.
- TowerSketch (TS-CMCU):
  - Extiende CountMin Sketch, organizando los contadores en “niveles” o torres de diferente tamaño.
  - Aplica conservative update, es decir, al insertar un k-mer solo se actualizan los contadores mínimos, evitando inflar artificialmente las frecuencias.
  - Maneja mejor los heavy hitters y distribuye el error de manera controlada.

## Precisión y error

- CS: ofrece estimaciones simétricas, lo que reduce el error relativo, pero puede generar más variabilidad en los resultados.
- TS-CMCU: tiende a mejorar la precisión práctica gracias al conservative update, reduciendo falsos positivos en los heavy hitters.

## Uso de memoria

- CS: Requiere un tamaño de sketch definido por parámetros (d, w). El error decrece con más memoria, pero no diferencia la importancia entre elementos.
- TS-CMCU: Usa memoria de forma más eficiente porque concentra contadores en distintos niveles, logrando un mejor balance entre exactitud y espacio.

## Ventajas y desventajas

- CS:
  - Buen rendimiento teórico para distribuciones con gran dispersión.
  - Simétrico: maneja bien valores positivos y negativos.

- Puede tener alta varianza en la estimación.
- TS-CMCU:
  - Más preciso en la práctica para heavy hitters.
  - Conservative update evita sobreestimación.
  - Mayor complejidad de implementación y ajuste de parámetros.

#### Aplicación en k-mers

- CS: adecuado cuando se requiere una primera aproximación rápida al conteo de k-mers frecuentes en secuencias genómicas.
- TS-CMCU: preferible cuando se busca un equilibrio más fuerte entre bajo error y uso razonable de memoria, especialmente útil en conjuntos de datos grandes y ruidosos.

## Conclusión

Los resultados obtenidos muestran que el algoritmo logra identificar un número considerable de heavy hitters en los archivos genómicos, aunque con un nivel variable de exactitud. Se observa que, en la mayoría de los casos, el recall es cercano a 1, lo que indica que casi todos los heavy hitters reales son detectados. Sin embargo, la precisión es significativamente menor, reflejando que junto a los verdaderos positivos aparecen también muchos falsos positivos. Esto provoca que el valor F1, que equilibra precisión y recall, sea generalmente bajo.

En términos prácticos, esto significa que el método es eficaz para no perder k-mers frecuentes (alta sensibilidad), pero menos fiable para discriminar cuáles son realmente representativos (baja especificidad). Por lo tanto, se requiere una calibración más fina de los parámetros del sketch y posiblemente la combinación con otras técnicas, a fin de reducir falsos positivos y mejorar la precisión sin sacrificar el recall.