

Ecosistema Hadoop

Máster en Data Science y Big Data

Jorge López-Malla Matute

jlmalla@geoblink.com

Febrero 2022

Presentación



- Jorge López-Malla Matute
- Puesto actual:
 - Senior Data Engineer en Geoblink
- Experiencia docente:
 - Profesor en diversos Masters de Big Data durante los últimos 6 años
 - Profesor de Tecnologías Masivas en ICAI
- Años trabajados en Big Data: 9 años

Índice

Contenido

1. Componentes

- Apache Hive
- Apache HBase

Componentes

Apache Hive

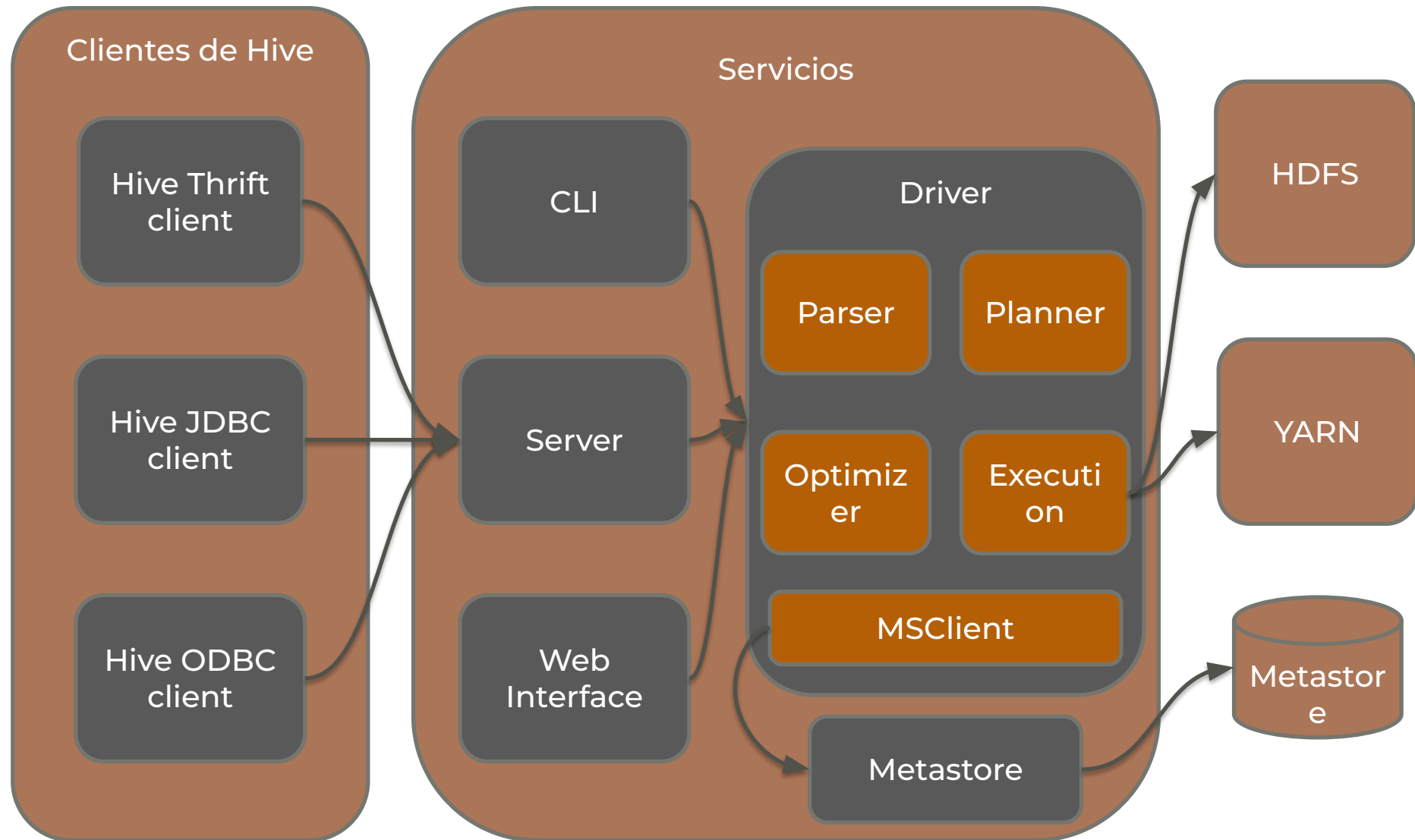
Hive- Introducción

- En el inicio **Map & Reduce** sólo estaba enfocado al desarrollo
- Con el tiempo más casos de uso se fueron desarrollando y se vio la necesidad de lenguajes de scripting
- El primer proyecto de Hadoop en scripting fue **Apache Pig** (2008)
- Los usuarios tenían que seguir aprendiendo otro lenguaje
- No podía conectarse con otros servicios externos con facilidad como lo hubiera hecho el SQL

Hive- Introducción

- En 2010 Facebook libera un proyecto que convierte sentencias SQL en trabajos de Map&Reduce
- Se libera Apache Hive y se extiende su uso
- Al ser un intérprete de SQL permite que herramientas que ya usaban SQL pudieran procesar cantidades masivas de datos
- Hive además permite la compartición de tablas de datos entre distintos procesos de una manera sencilla
- Los proyectos posteriores de procesamiento distribuido toman Hive como referencia

Hive- Arquitectura



Hive- Funcionamiento

- **Hive** interpreta sentencias SQL en un motor de procesamiento distribuido
- **Hive** necesita de una base de datos de metadatos
- En ella se guardan los metadatos de las tablas
- Las tablas tienen que almacenarse en un almacenamiento compatible con el motor de procesamiento
- **Hive** interpreta la secuencia y la convierte en trabajo distribuido más eficiente

Hive- Funcionamiento

Origen, Destino, Fecha, Pasajeros

NYC,SFO,20200101,120

SFO,LAX,20200101,55

NYC,BOS,20200203,100

MAD,BCN,20200301,120

NYC,SFO,20200203,110

NYC,MAD,20200301,200

NYC,BCN,20200301,220

BOS,LAS,20200308,115

MAD,NYC,20200308,215

LAS,MAD,20200601,55

BCN,NYC,20200605,60

LAS,SFO,20200607,110

BCN,MAD,20200705,100

BOS,NYC,20200708,108

Map1

Entrada:

Vuelos en csv

Proceso:

Partir la línea por ,

Salida:

clave: origen y destino

valor: 1

Reduce1

Entrada:

((Origen, Destino), 1)

Proceso:

Sumar los 1's.

Salida

clave: (Origen, Destino)

valor: contador de

vuelos

Map2

Entrada:

((Origen, Destino), contador)

Proceso:

Filtrar los vuelos con destino Boston

Salida

clave:

(Origen, Destino)

valor: contador

Reduce2

Entrada:

((Origen, Destino), contador)

Proceso:

Ninguno

Salida

clave: Origen

valor: contador

Map3

Entrada:

(Origen, contador)

Proceso:

Ninguno

Salida

clave: null

valor: (Origen,

contador)

Reduce3

Entrada:

(null, (Origen, contador))

Proceso:

Seleccionar el Origen con mayor numero de vuelos

Salida

clave: Origen

valor: contador

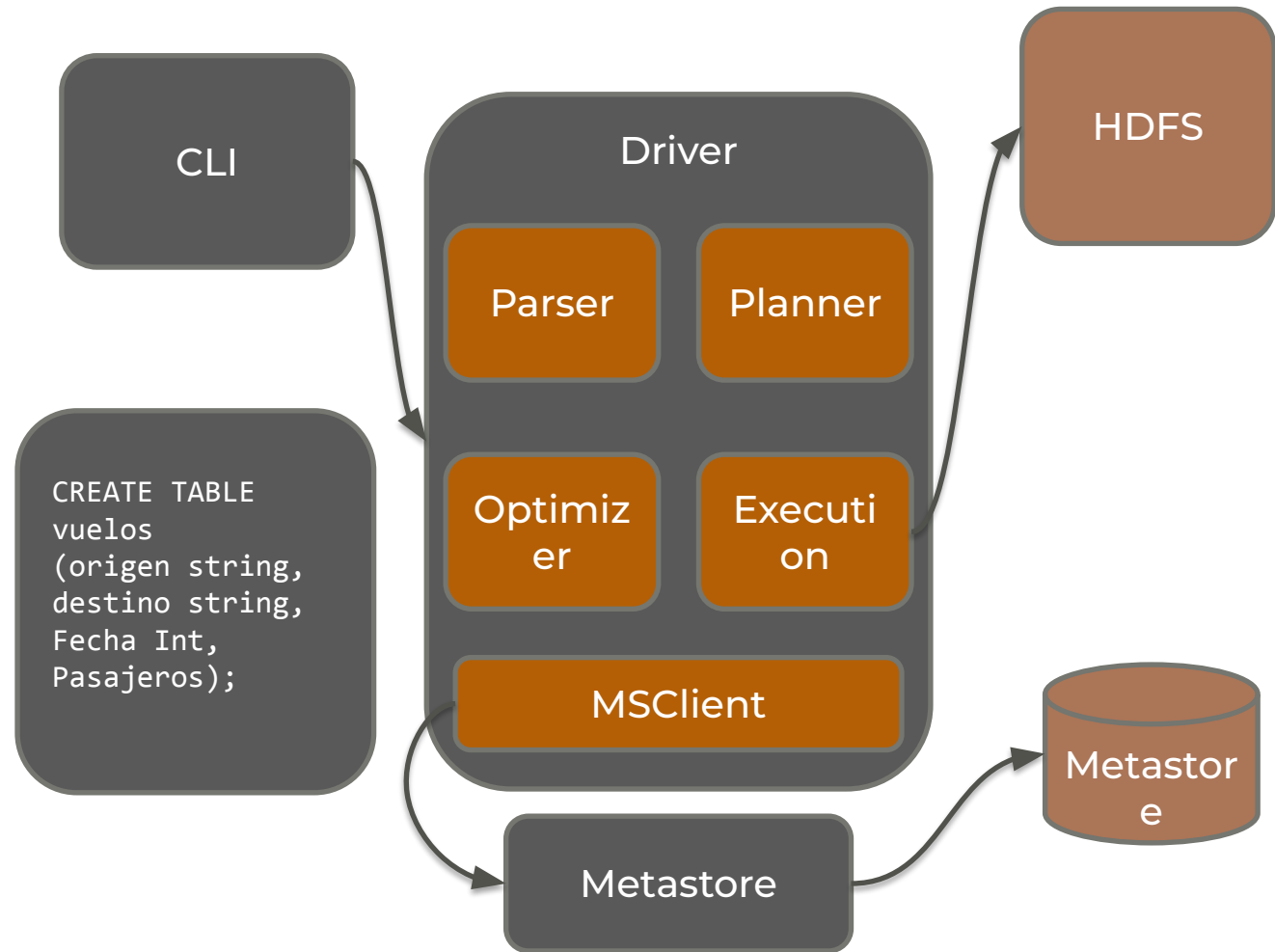
maximo

Hive- Funcionamiento

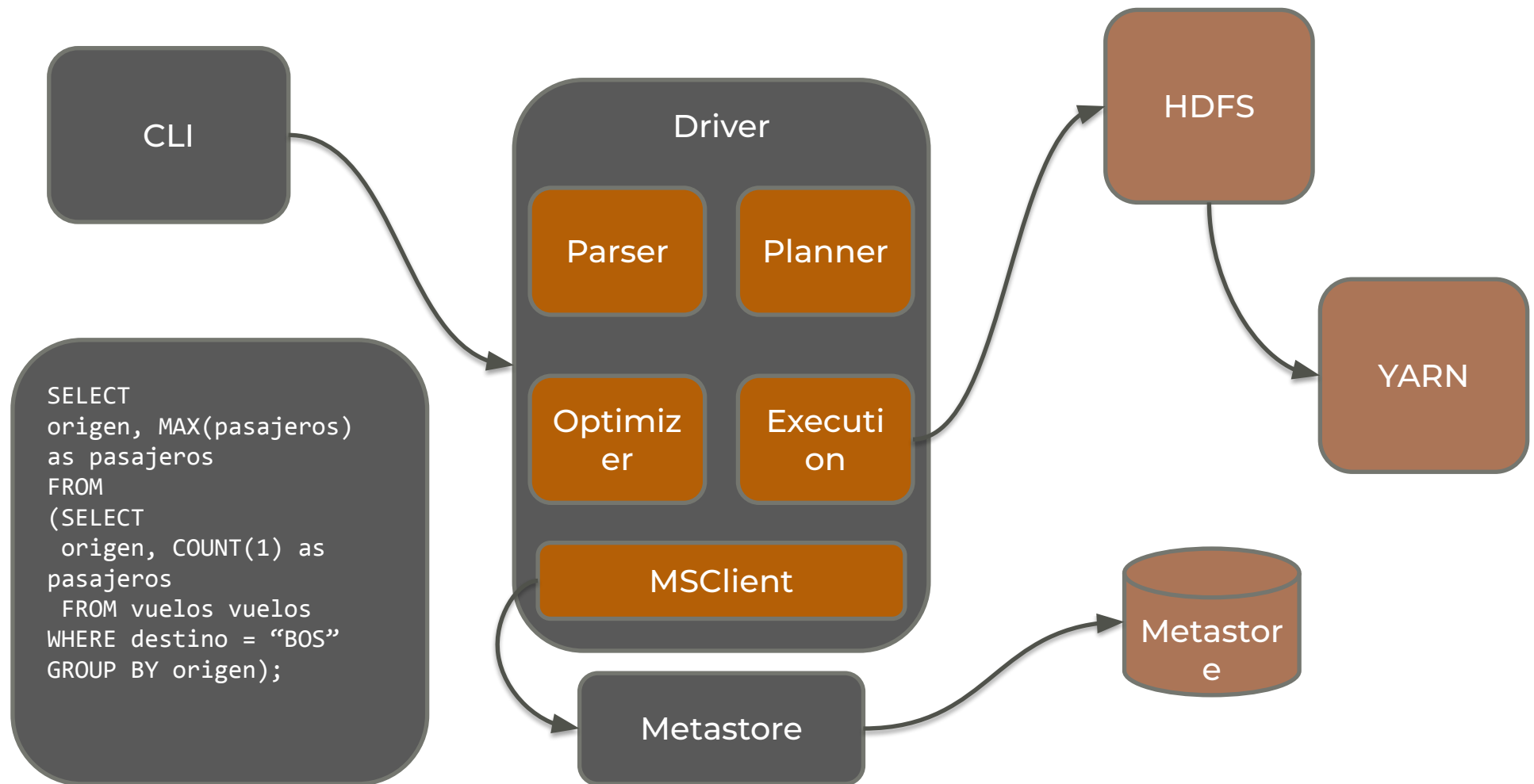
```
Origen, Destino, Fecha, Pasajeros  
NYC, SFO, 20200101, 120  
SFO, LAX, 20200101, 55  
NYC, BOS, 20200203, 100  
MAD, BCN, 20200301, 120
```

```
NYC, SFO, 20200203, 110  
NYC, MAD, 20200301, 200  
NYC, BCN, 20200301, 220  
BOS, LAS, 20200308, 115  
MAD, NYC, 20200308, 215
```

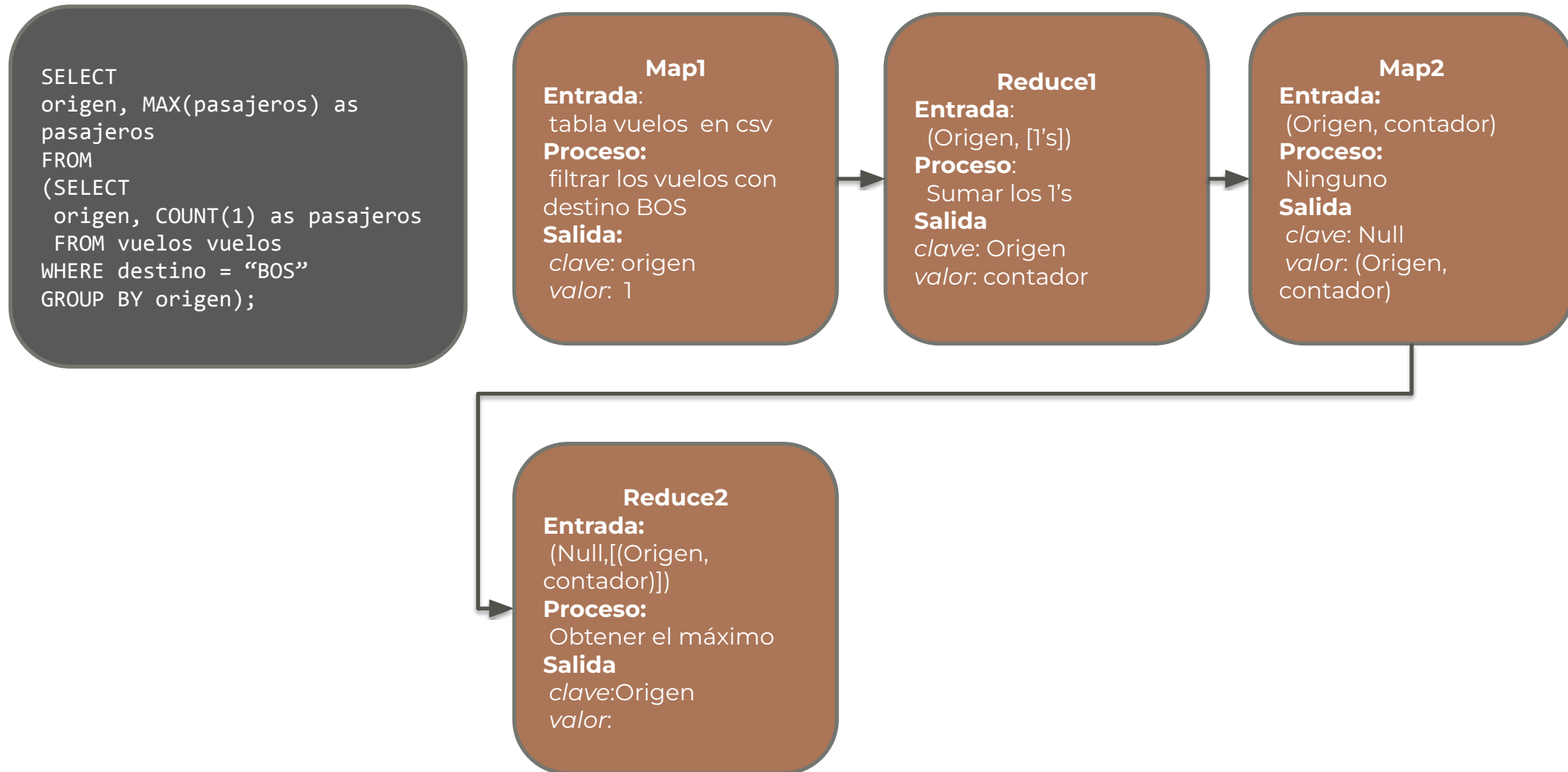
```
LAS, MAD, 20200601, 55  
BCN, NYC, 20200605, 60  
LAS, SFO, 20200607, 110  
BCN, MAD, 20200705, 100  
BOS, NYC, 20200708, 108
```



Hive- Funcionamiento



Hive- Funcionamiento



Hive- Introducción

- La solución de **Hive** da la solución óptima mirando el planificador
- El planificador de queries no entiende de datos ni de escalabilidad
- Puede hacer algún tipo de optimización usando el almacenamiento
- ¿Y si los vuelos con destino a Boston suman tanto como para llenar la partición más grande de nuestro sistema?
- Se podría solucionar con tablas intermedias

Apache Hbase

HBase - Introducción

- Map & Reduce procesa grandes cantidades de datos en un tiempo aceptable
- Hive nos proporciona un lenguaje SQL
- Si juntamos las dos tecnologías, ¿tendríamos una Base de Datos relacional con todas ventajas del Big Data?
- La respuesta es NO
- Juntando una tecnología de cómputo que admita cantidades masivas de datos y un proyecto que nos permita hacer SQL **NO** nos permite hacer queries en tiempo *online* “sólo” facilita el procesamiento

HBase - Introducción

- A raíz de las nuevas formas de almacenar se empiezan a pensar nuevas formas de consultar esa información
- El principal problema que pretende resolver es el acceso aleatorio a unos datos concretos
- Para ello usa HDFS como sistema de ficheros subyacente y una estrategia **Columnar**
- Con ello se consigue un acceso “rápido” a datos aleatorios en colecciones masivas de información
- Posteriormente nacen tecnologías de almacenamiento basadas en su misma filosofía supliendo algunas de sus carencias (Apache Cassandra)

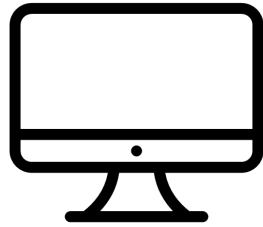
HBase - Introducción

- A raíz de las nuevas formas de almacenar se empiezan a pensar nuevas formas de consultar esa información
- El principal problema que pretende resolver es el acceso aleatorio a unos datos concretos
- Para ello usa HDFS como sistema de ficheros subyacente y una estrategia **Columnar**
- Con ello se consigue un acceso “rápido” a datos aleatorios en colecciones masivas de información
- Posteriormente nacen tecnologías de almacenamiento basadas en su misma filosofía supliendo algunas de sus carencias (Apache Cassandra)

HBase - Terminología

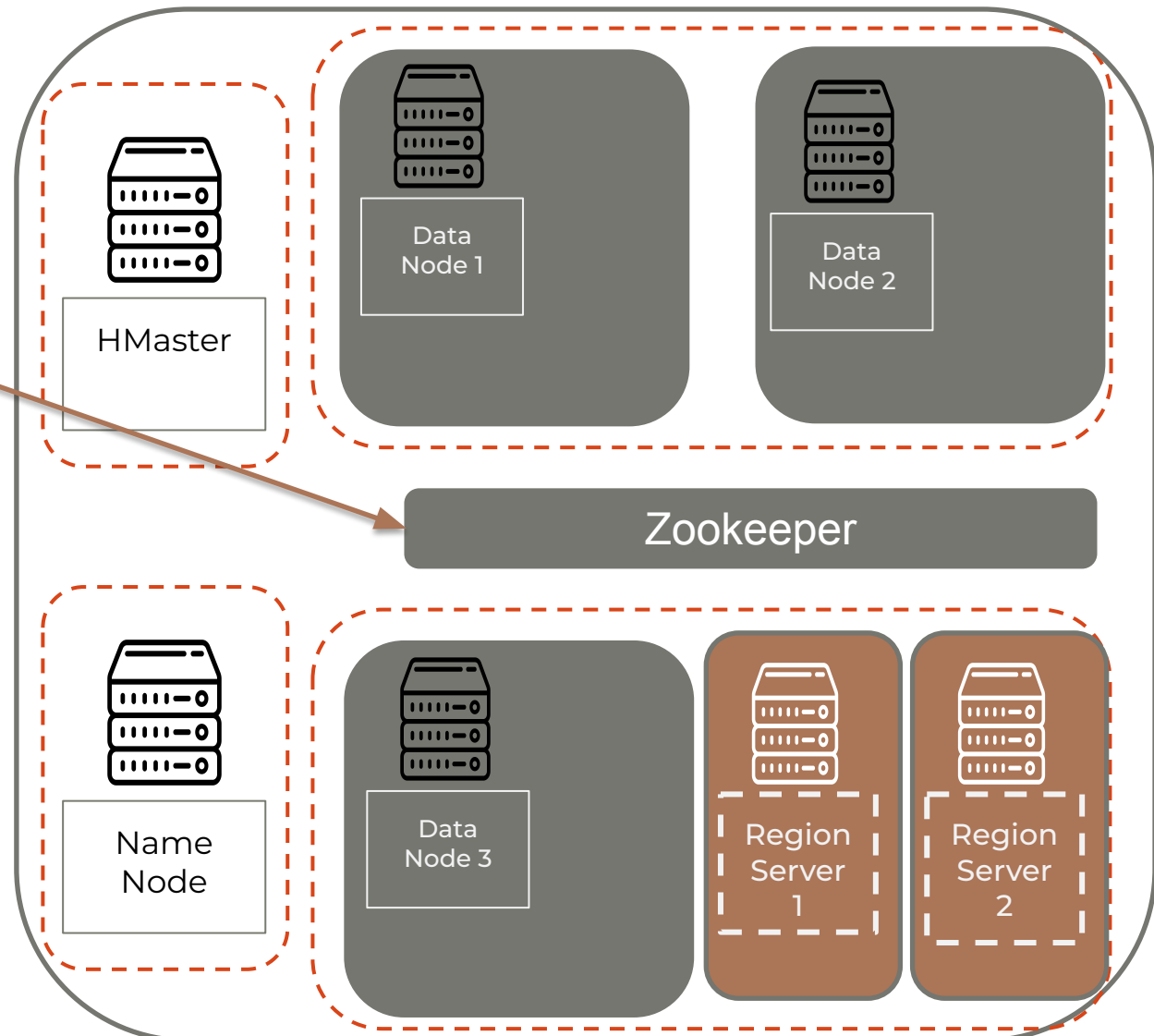
- Hbase tiene la siguiente terminología:
 - **Table:** Conjunto de **rows**
 - **Row:** Conjunto de datos agrupados en columnas. Se compone de una **rowKey** y varias **column families**
 - **RowKey:** Column de una **row** que hace única a la misma.
 - **ColumnFamily:** Serie de Columns de una row. Su estructura no se tiene que compartir entre distintas **rows** de una misma **table**

HBase Funcionamiento

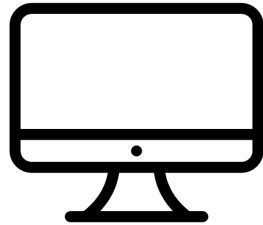


Client-1

```
hbase create table  
'electric_company_cups_regiones',  
'region_data', 'tarificacion_data'
```

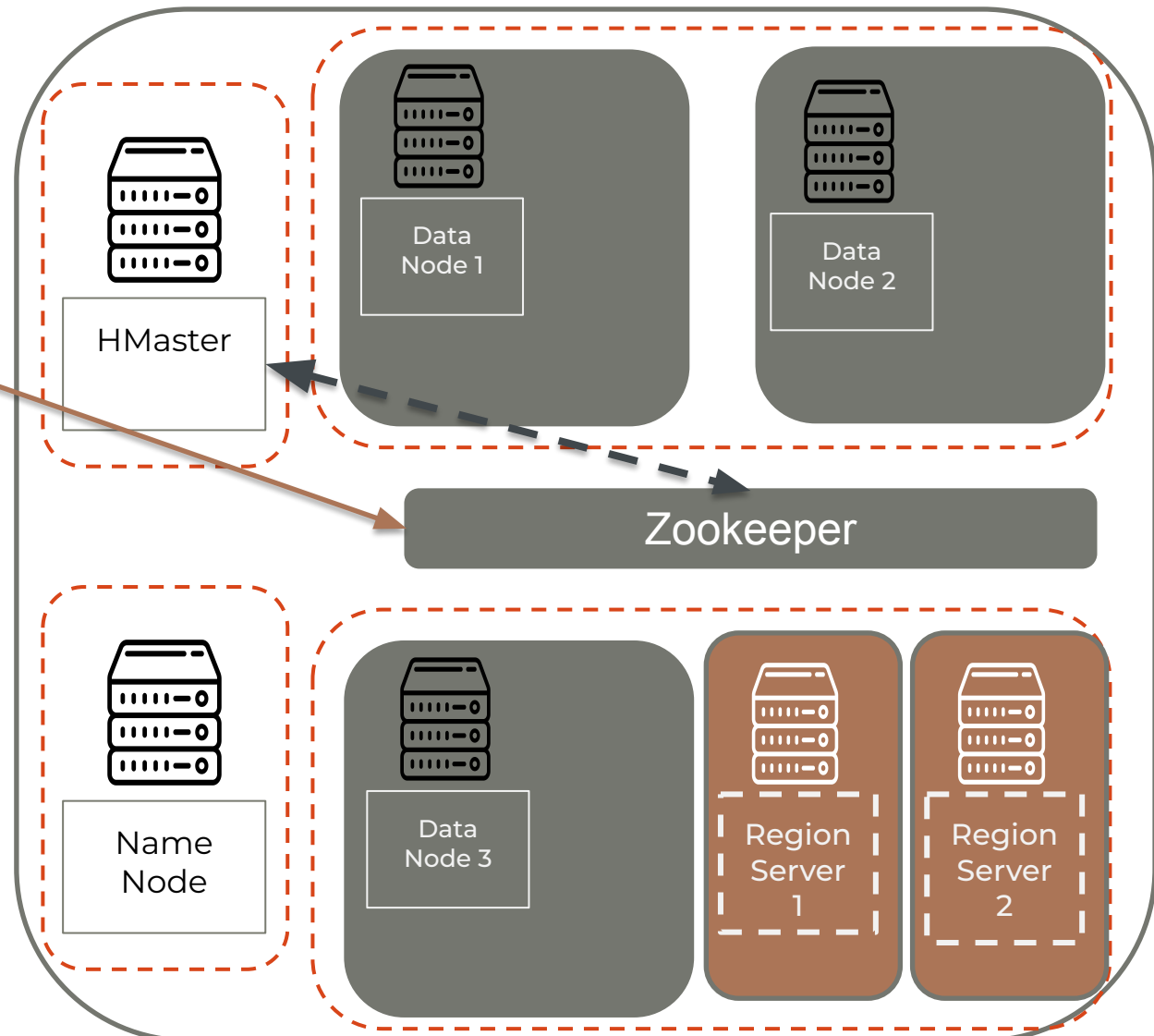


HBase Funcionamiento

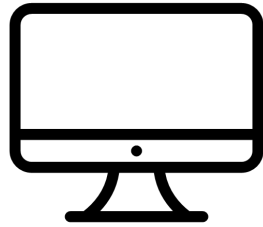


Client-1

```
hbase create table  
'electric_company_cups_regiones',  
'region_data', 'tarificacion_data'
```

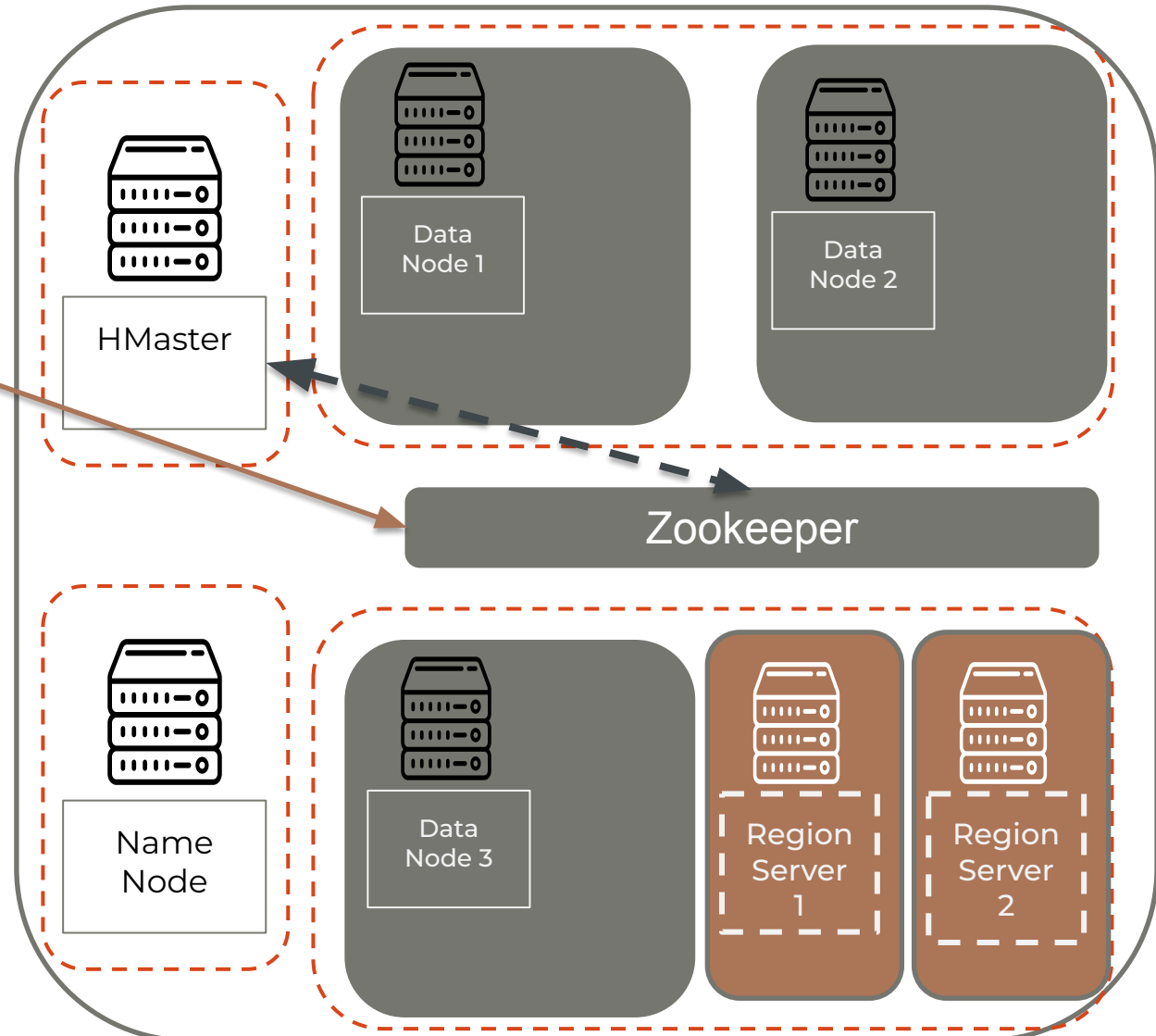


HBase Funcionamiento

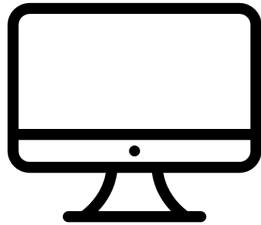


Client-1

```
hbase create table  
'electric_company_cups_regiones',  
'region_data', 'tarificacion_data'
```



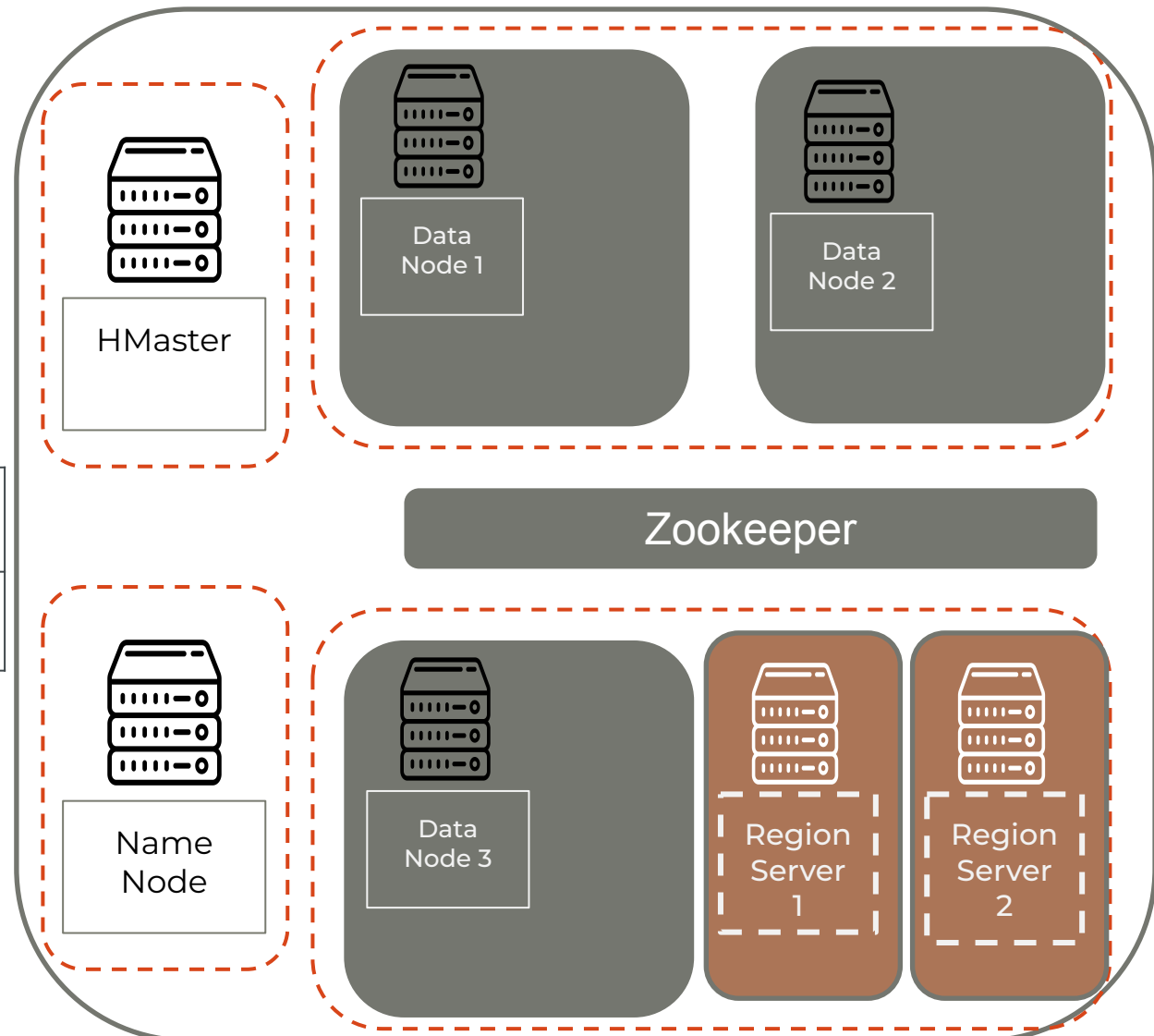
HBase Funcionamiento



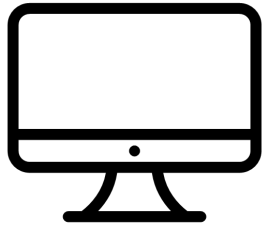
Client-1

```
hbase create table  
'electric_company_cups_regiones',  
'region_data', 'tarificacion_data'
```

Row	Region_data	Tarificacion_data	timestamp



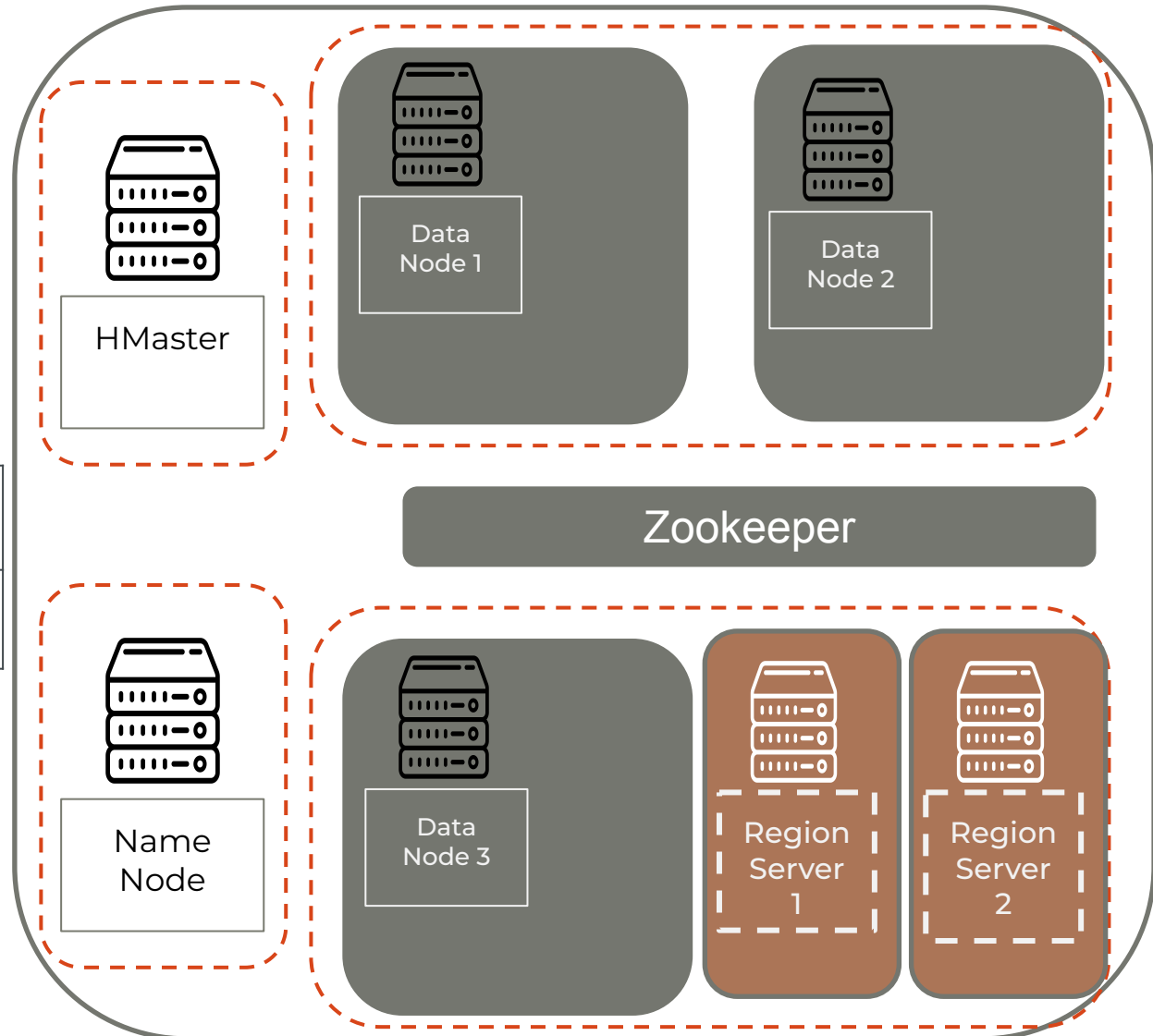
HBase Funcionamiento



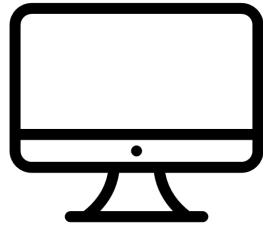
Client-1

```
hbase put  
electric_company_cups_regiones  
'0001_mun', 'region_data:valor_region'  
'Colmenar Viejo'
```

Row	Region_data	Tarificacion_data	timestamp



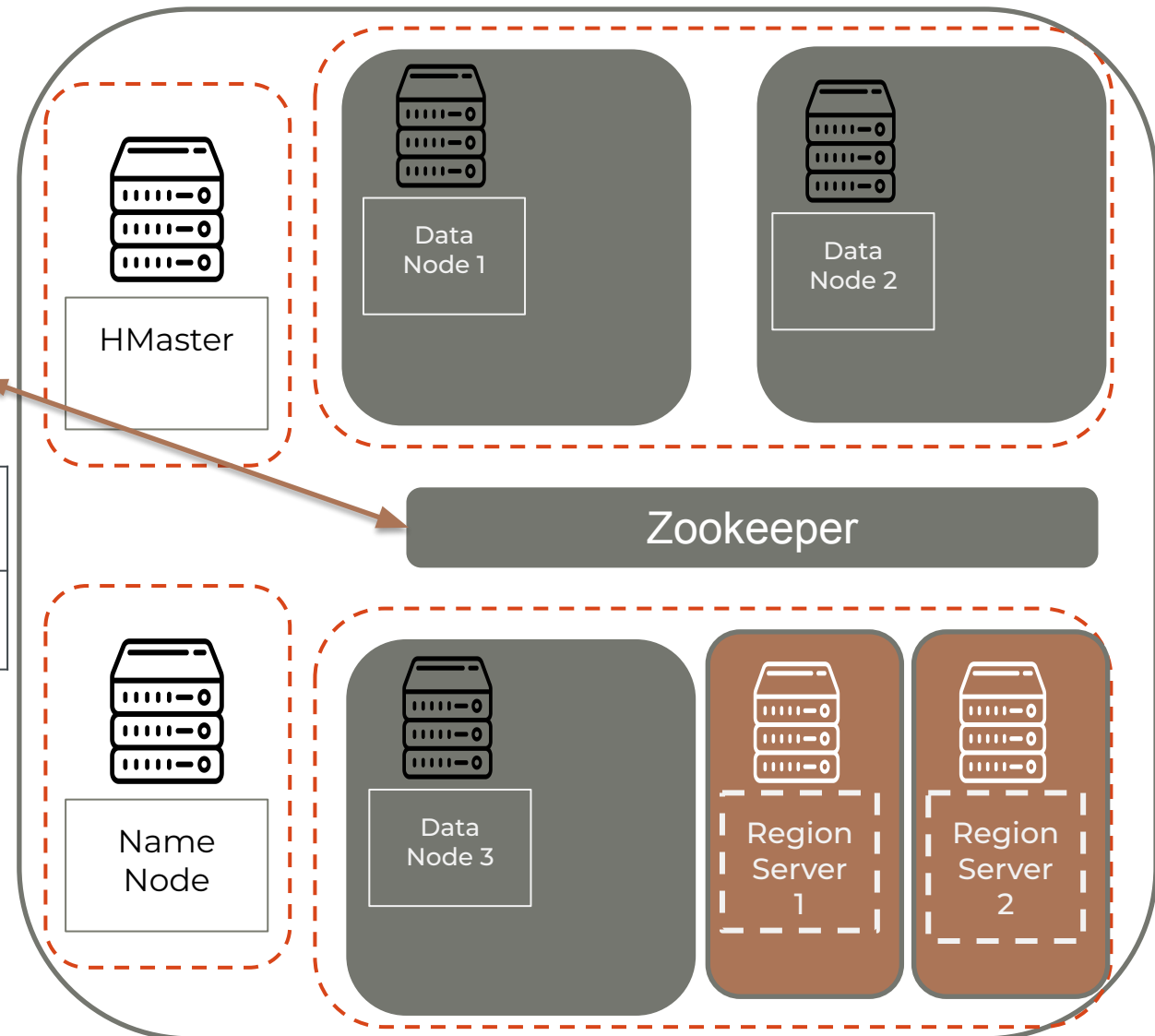
HBase Funcionamiento



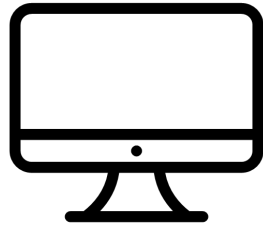
Client-1

```
hbase put  
electric_company_cups_regiones  
'0001_mun', 'region_data:valor_region'  
'Colmenar Viejo'
```

Row	Region_data	Tarificacion_data	timestamp



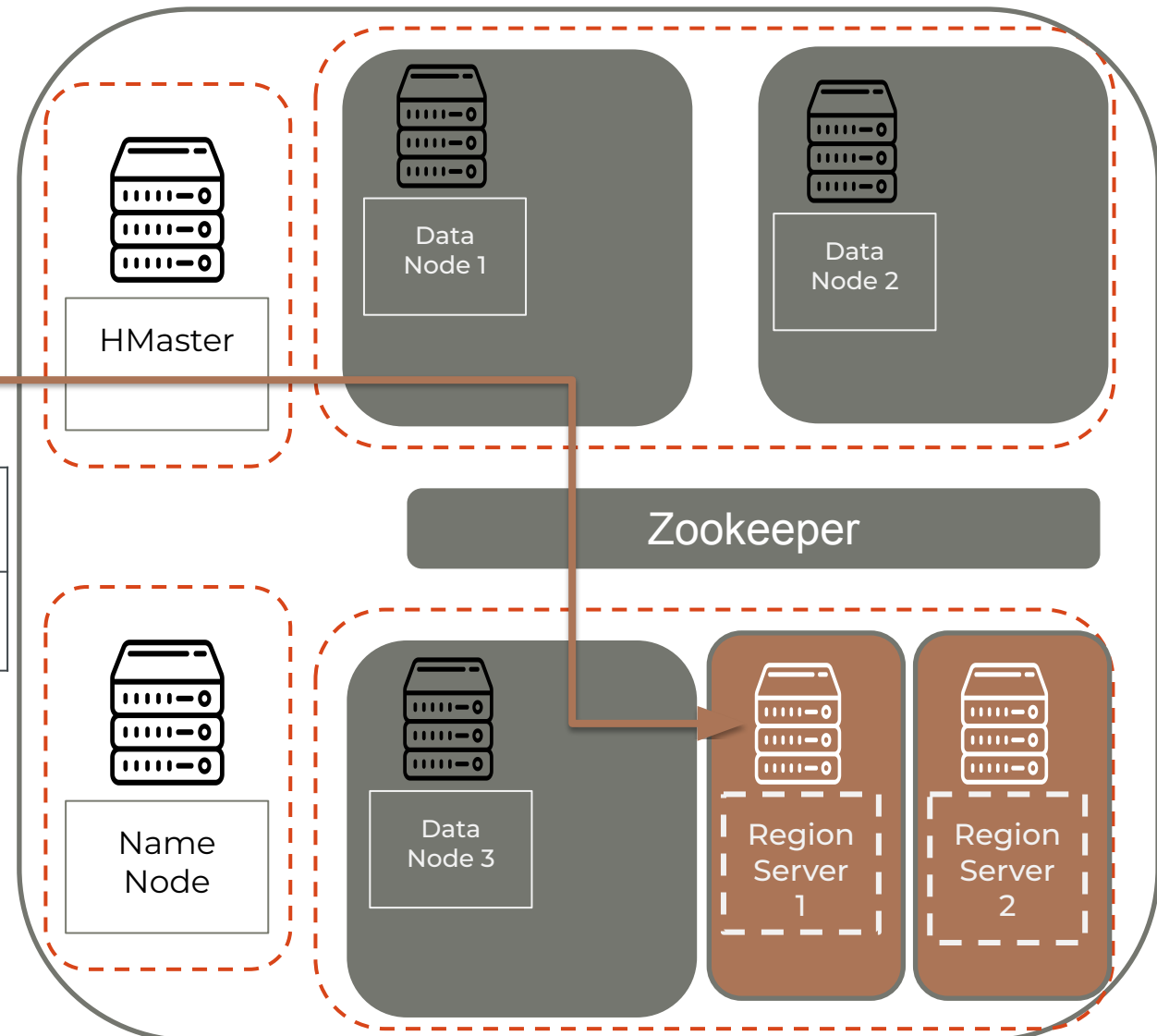
HBase Funcionamiento



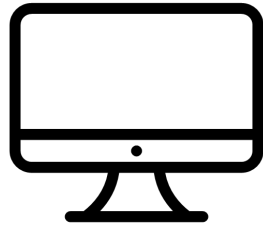
Client-1

```
hbase put  
electric_company_cups_regiones  
'0001_mun', 'region_data:valor_region'  
'Colmenar Viejo'
```

Row	Region_data	Tarificacion_data	timestamp



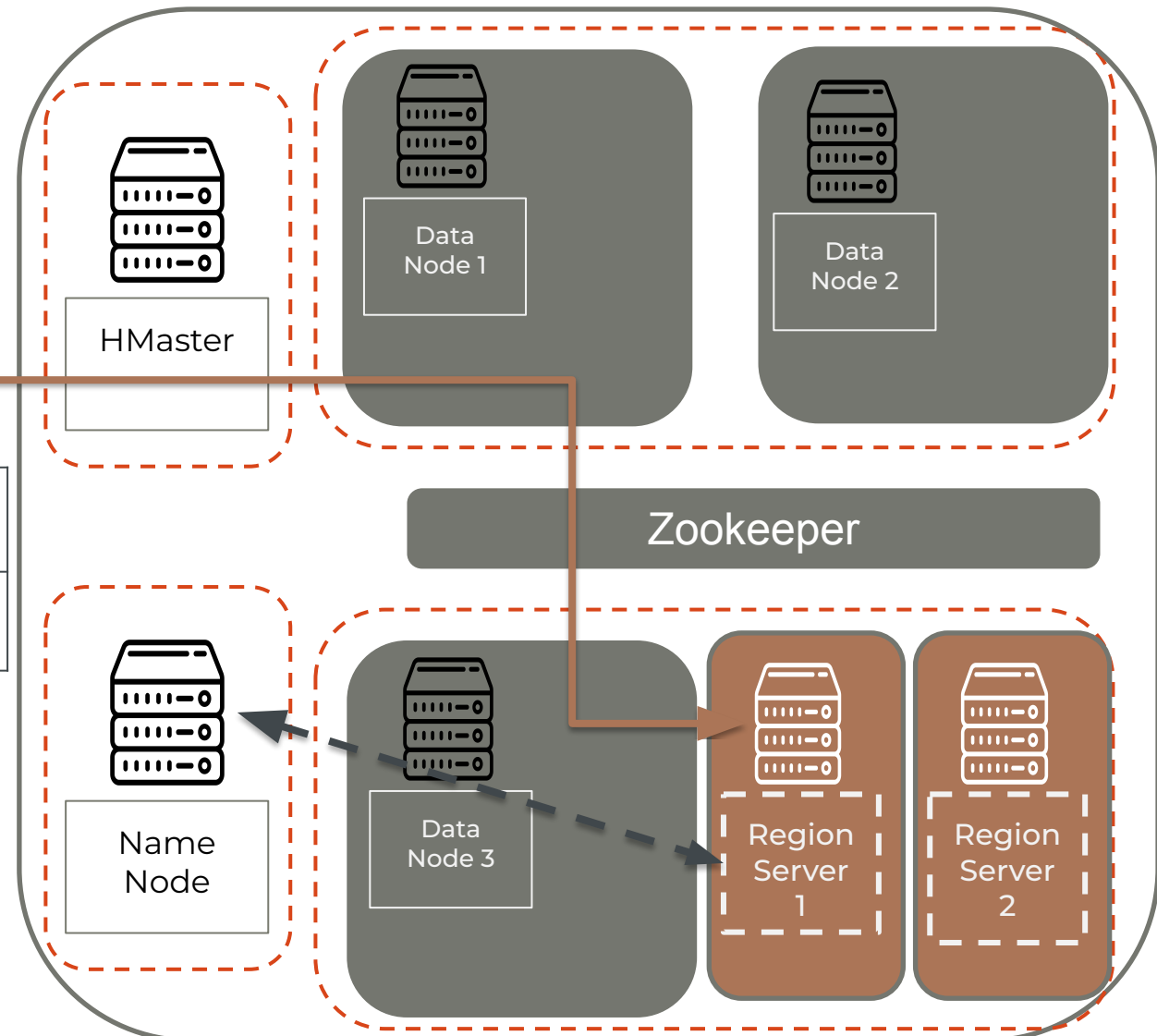
HBase Funcionamiento



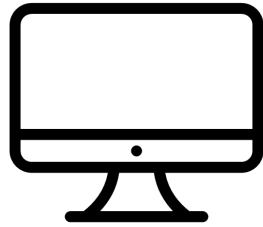
Client-1

```
hbase put  
electric_company_cups_regiones  
'0001_mun', 'region_data:valor_region'  
'Colmenar Viejo'
```

Row	Region_data	Tarificacion_data	timestamp



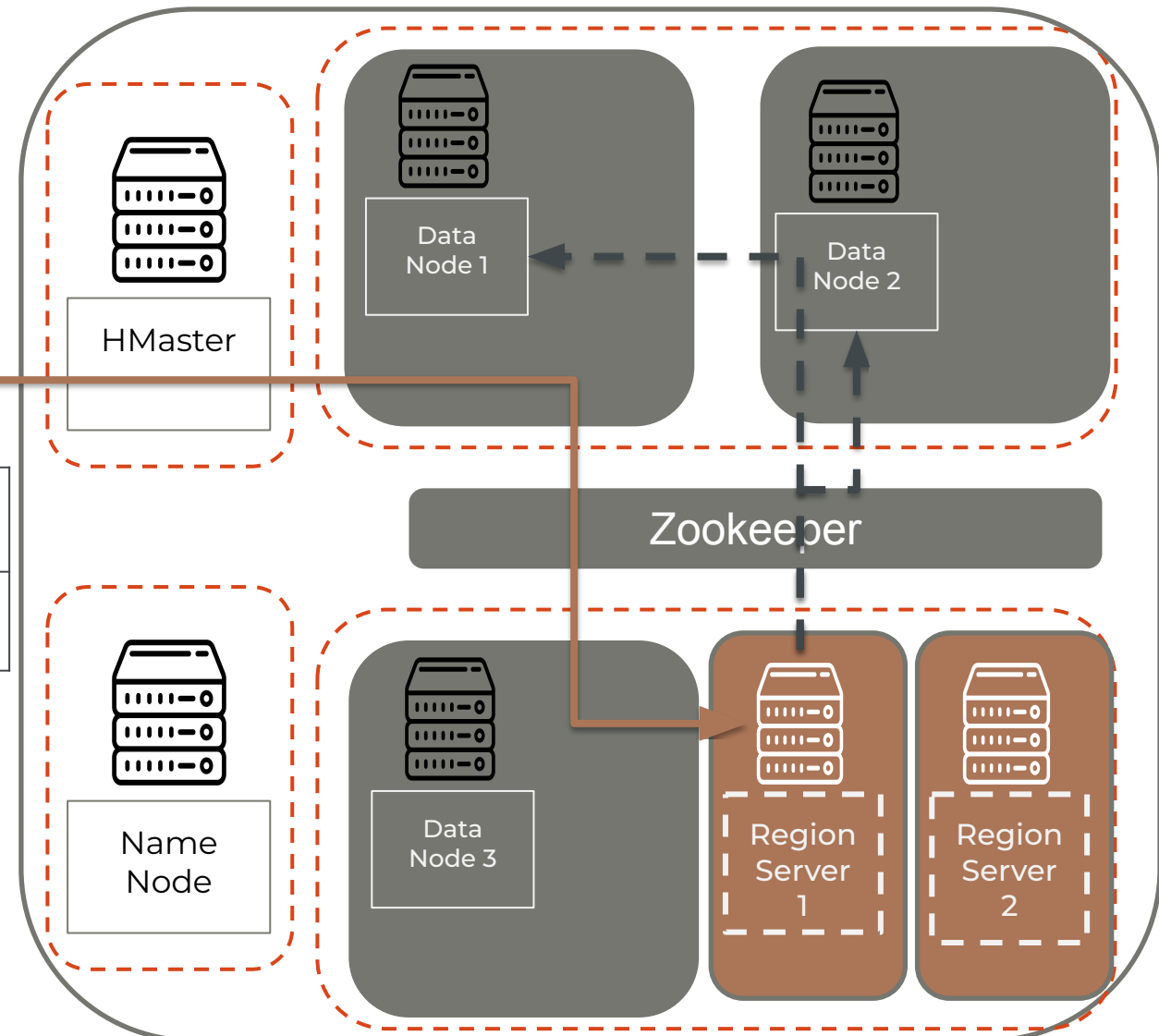
HBase Funcionamiento



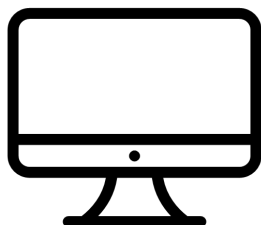
Client-1

```
hbase put  
electric_company_cups_regiones  
'0001_mun', 'region_data:valor_region'  
'Colmenar Viejo'
```

Row	Region_data	Tarificacion_data	timestamp



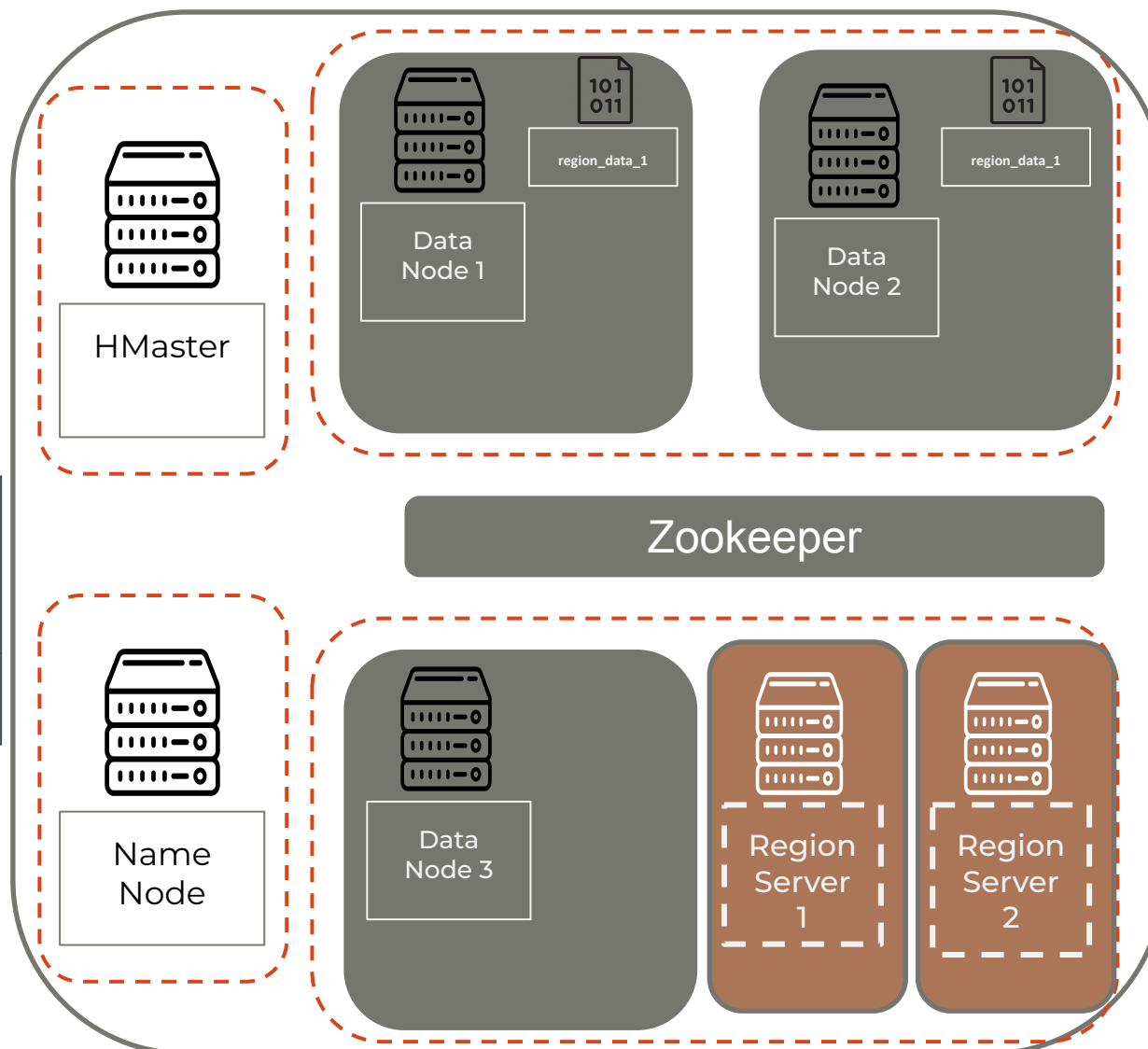
HBase Funcionamiento



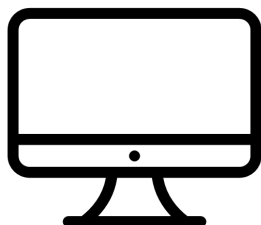
Client-1

```
hbase put
electric_company_cups_regiones
'0001_mun', 'region_data:valor_region'
'Colmenar Viejo'
```

Row	Region_data	Tarificacio n_data	timestamp
	valor_region		
0001_mun	Colmenar Viejo		1631952266000



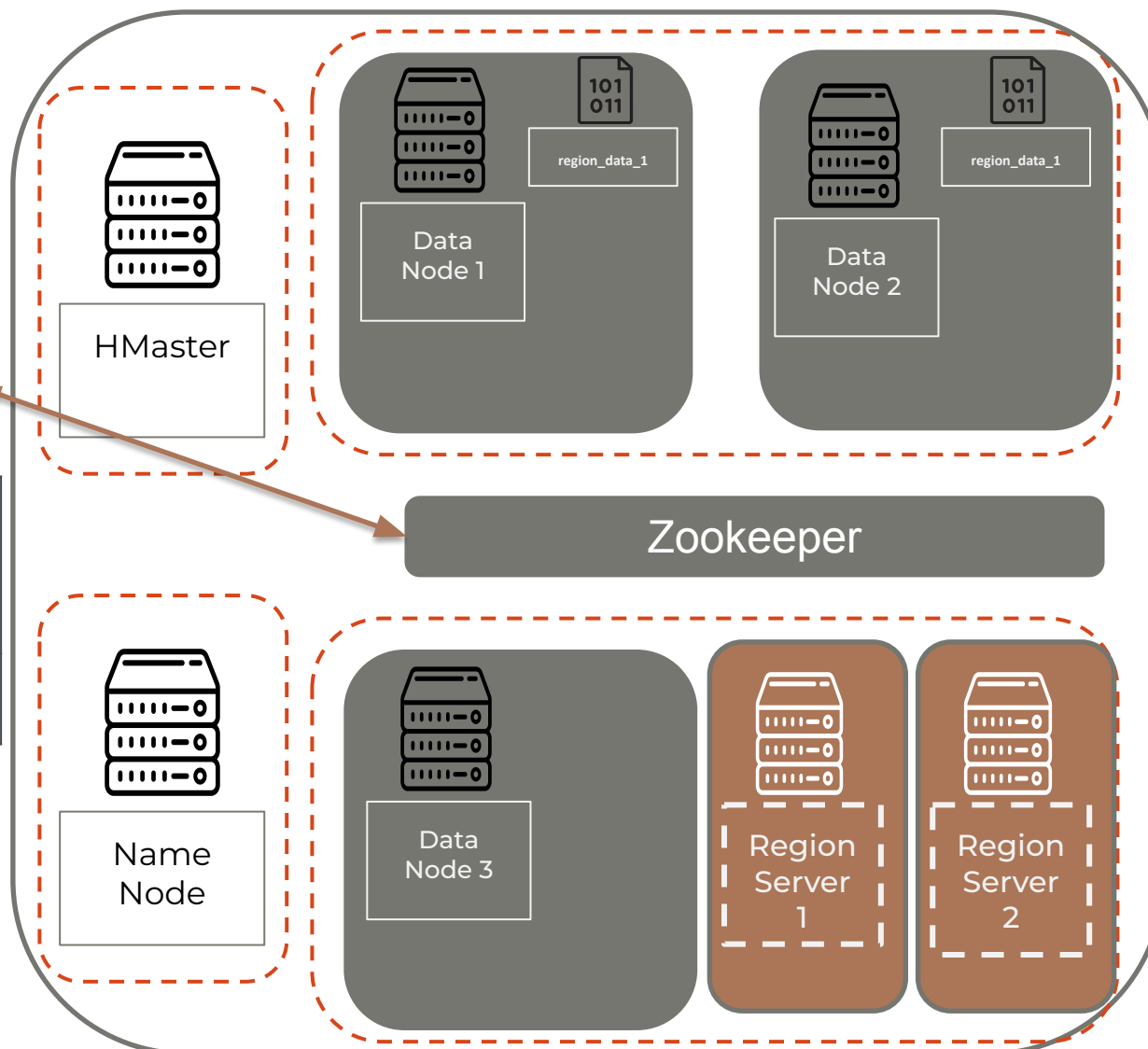
HBase Funcionamiento



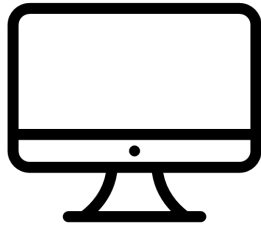
Client-1

```
hbase put
electric_company_cups_regiones
'0001_mun', 'Tarificacion_data':'valle
10.34'
```

Row	Region_data	Tarificacio n_data	timestamp
	valor_region		
0001_mun	Colmenar Viejo		1631952266000



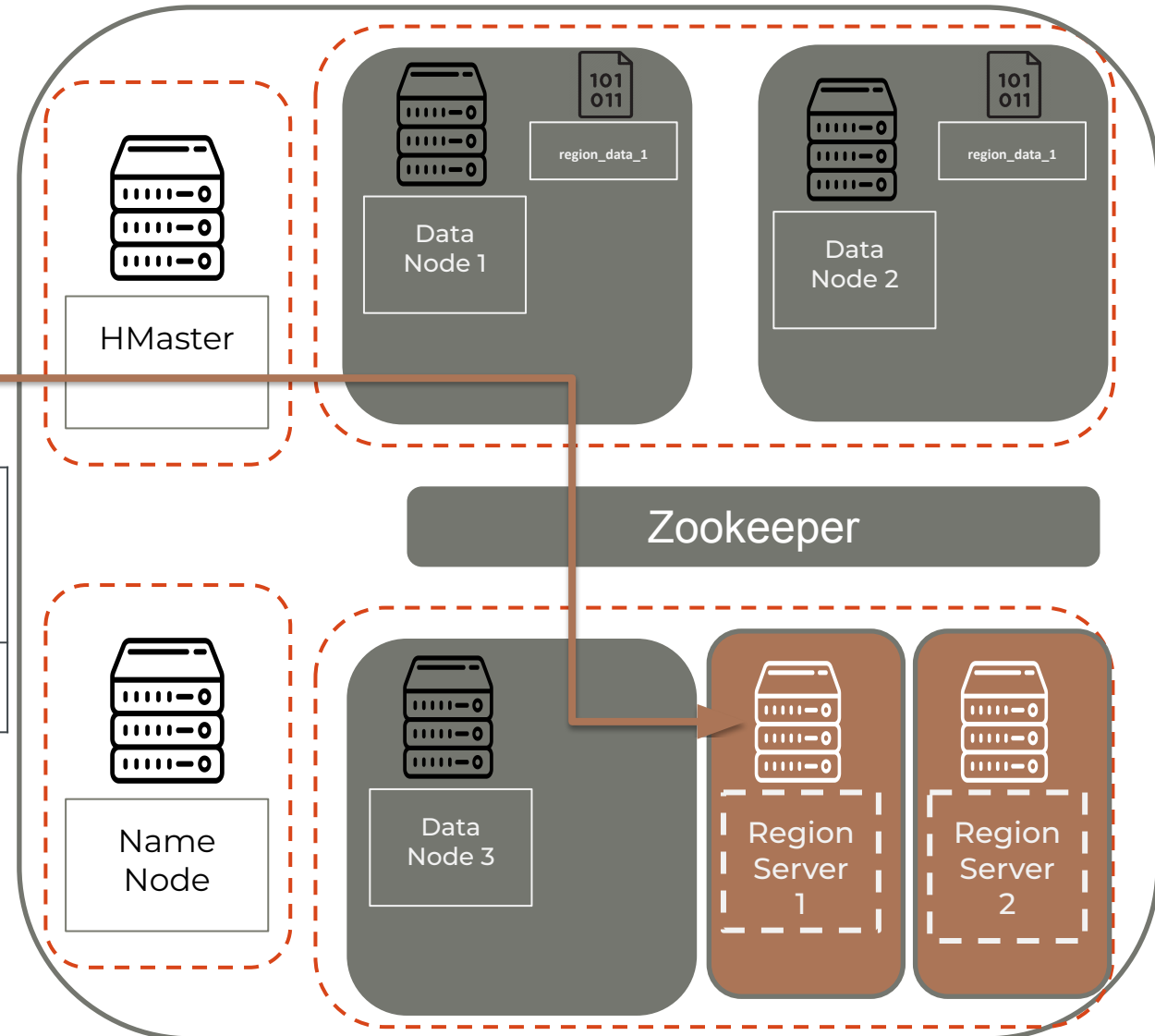
HBase Funcionamiento



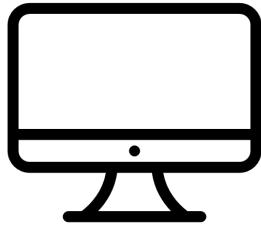
Client-1

```
hbase put
electric_company_cups_regiones
'0001_mun', 'Tarificacion_data': 'valle
10.34'
```

Row	Region_data	Tarificacio n_data	timestamp
	valor_region		
0001_mun	Colmenar Viejo		1631952266000



HBase Funcionamiento



Client-1

```
hbase put
electric_company_cups_regiones
'0001_mun', 'Tarificacion_data': 'valle
10.34'
```

Row	Region_data	Tarificacio n_data	timestamp
	valor_region		
0001_mun	Colmenar Viejo		1631952266000

HMaster

Data
Node 1

region_data_1

101
011

Data
Node 2

region_data_1

101
011

Zookeeper

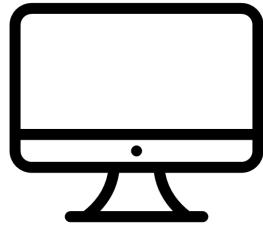
Name	Node
1	1
2	2
3	3
4	4
5	5
6	6
7	7
8	8
9	9
10	10
11	11
12	12
13	13
14	14
15	15
16	16
17	17
18	18
19	19
20	20
21	21
22	22
23	23
24	24
25	25
26	26
27	27
28	28
29	29
30	30
31	31
32	32
33	33
34	34
35	35
36	36
37	37
38	38
39	39
40	40
41	41
42	42
43	43
44	44
45	45
46	46
47	47
48	48
49	49
50	50
51	51
52	52
53	53
54	54
55	55
56	56
57	57
58	58
59	59
60	60
61	61
62	62
63	63
64	64
65	65
66	66
67	67
68	68
69	69
70	70
71	71
72	72
73	73
74	74
75	75
76	76
77	77
78	78
79	79
80	80
81	81
82	82
83	83
84	84
85	85
86	86
87	87
88	88
89	89
90	90
91	91
92	92
93	93
94	94
95	95
96	96
97	97
98	98
99	99
100	100

Data
Node 3

Region
Server
1

Region
Server
2

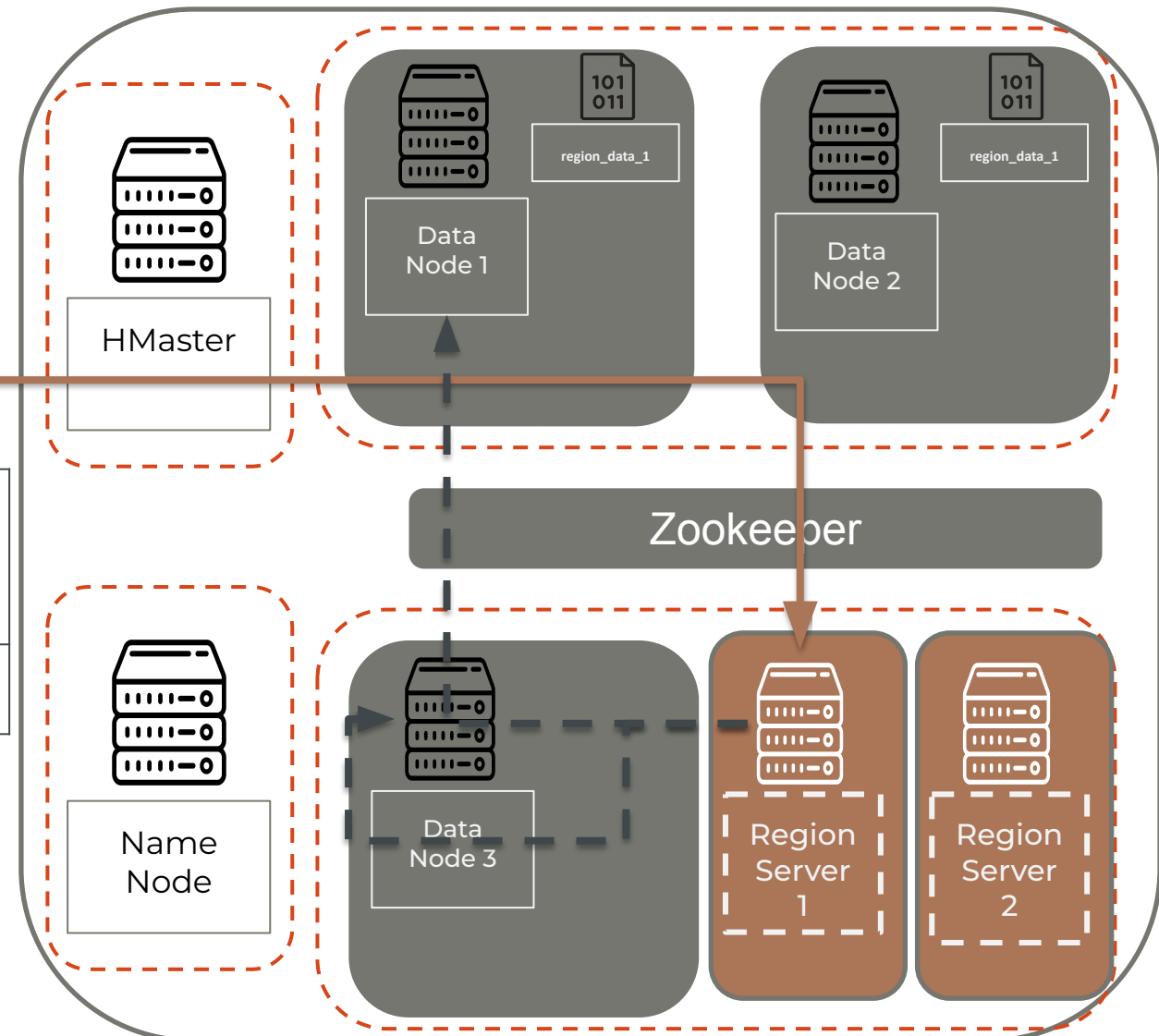
HBase Funcionamiento



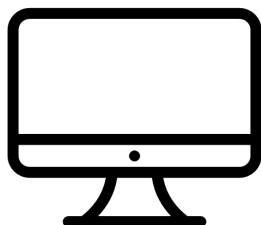
Client-1

```
hbase put
electric_company_cups_regiones
'0001_mun', 'Tarificacion_data':'valle
10.34'
```

Row	Region_data	Tarificacion_data	timestamp
	valor_region		
0001_mun	Colmenar Viejo		1631952266000



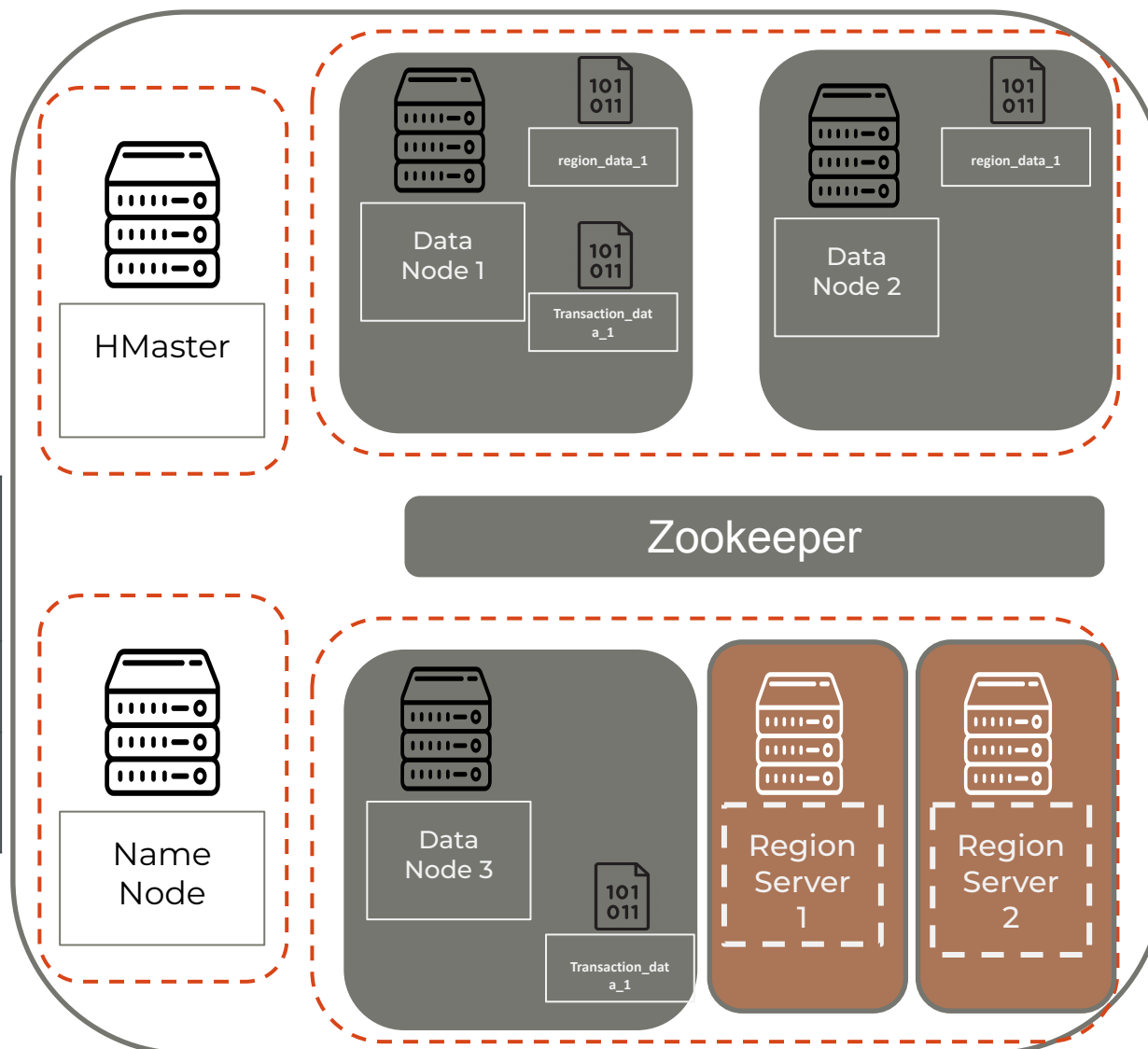
HBase Funcionamiento



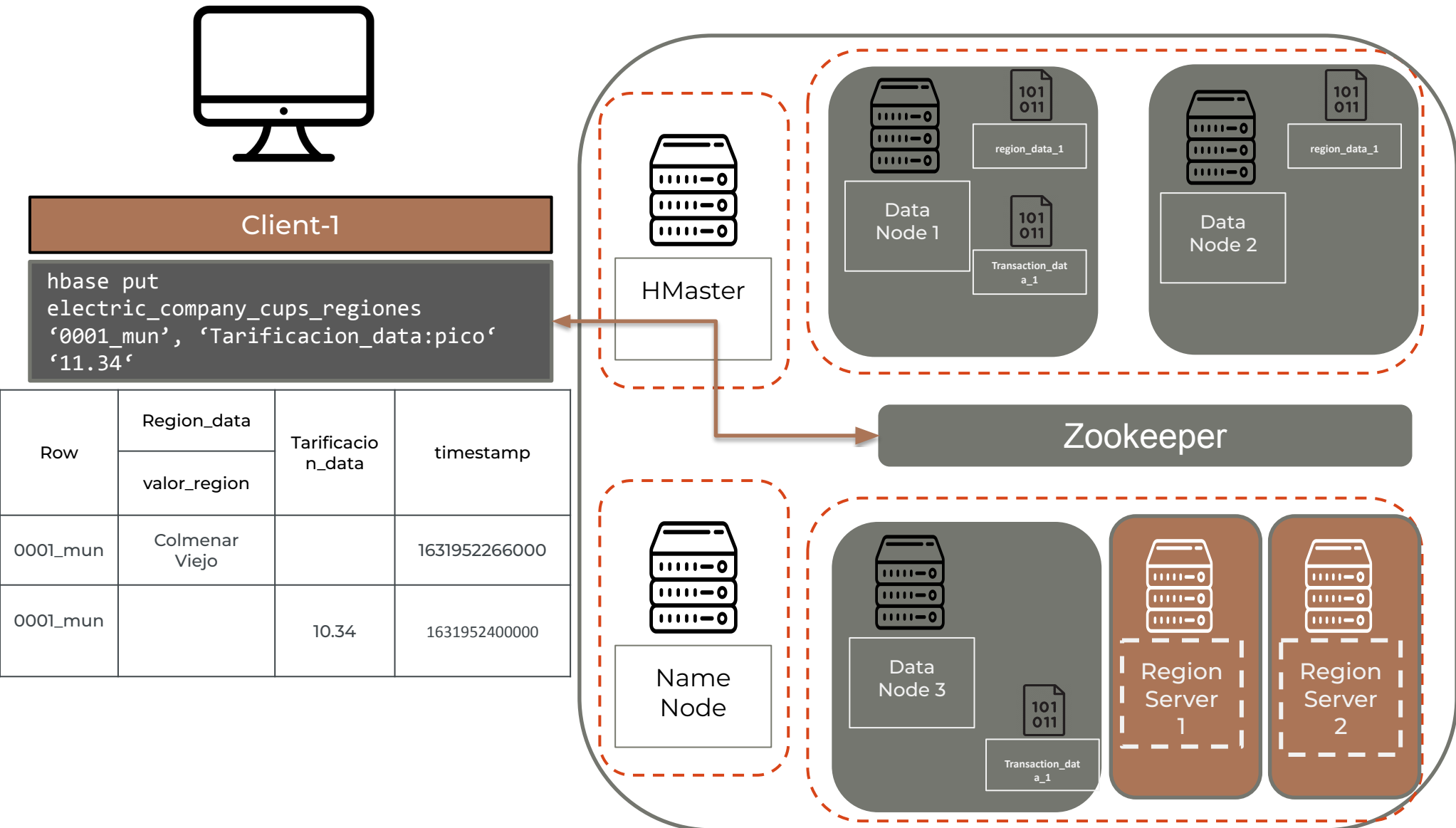
Client-1

```
hbase put
electric_company_cups_regiones
'0001_mun', 'Tarificacion_data':'valle
10.34'
```

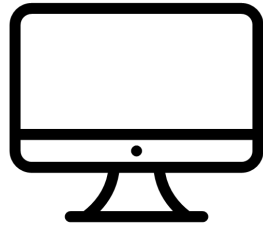
Row	Region_data	Tarificacio n_data	timestamp
	valor_region		
0001_mun	Colmenar Viejo		1631952266000
0001_mun		10.34	1631952400000



HBase Funcionamiento



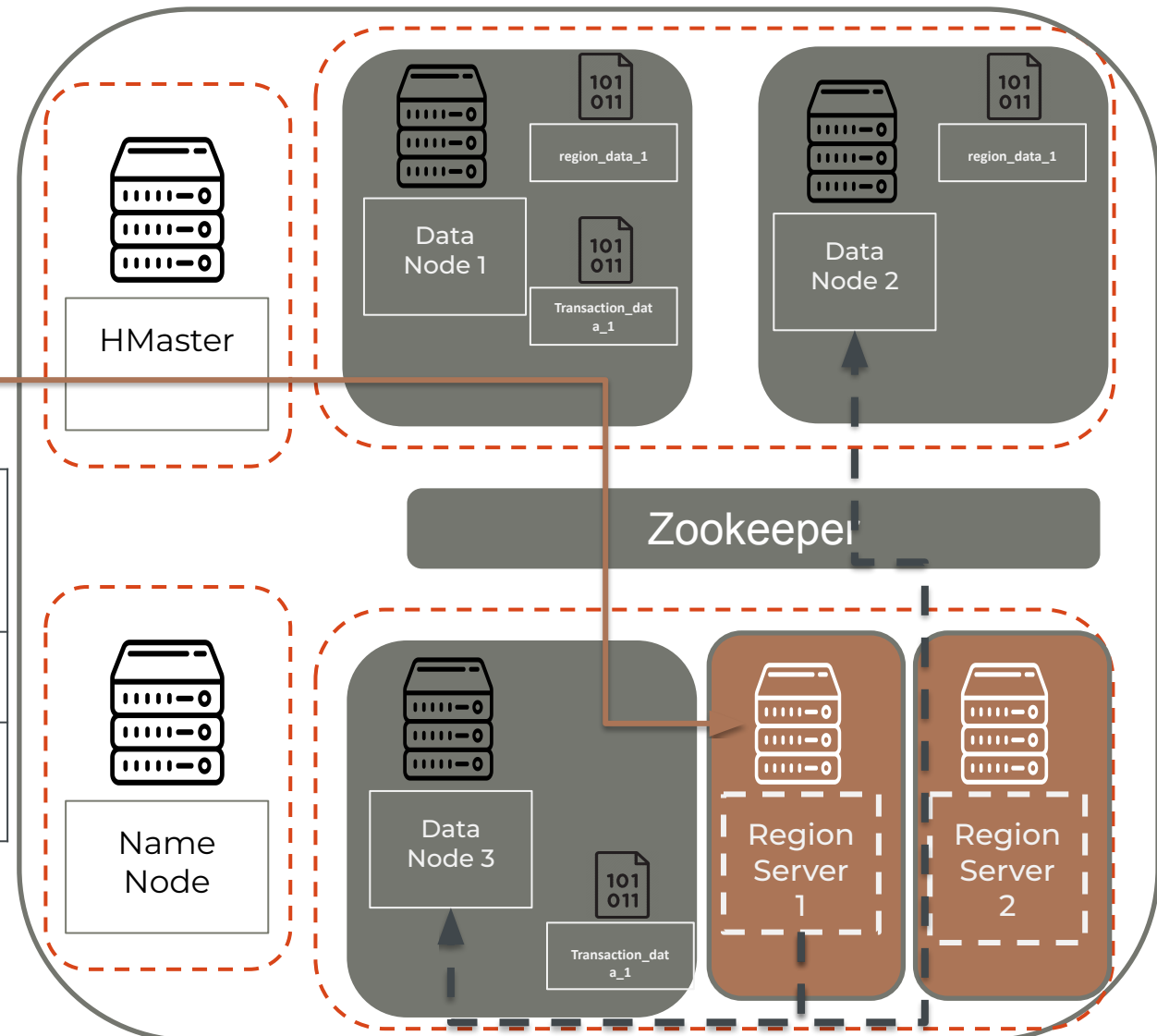
HBase Funcionamiento



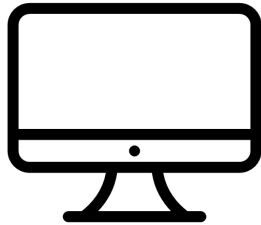
Client-1

```
hbase put
electric_company_cups_regiones
'0001_mun', 'Tarificacion_data:pico'
'11.34'
```

Row	Region_data	Tarificacio n_data	timestamp
	valor_region		
0001_mun	Colmenar Viejo		1631952266000
0001_mun		10.34	1631952400000



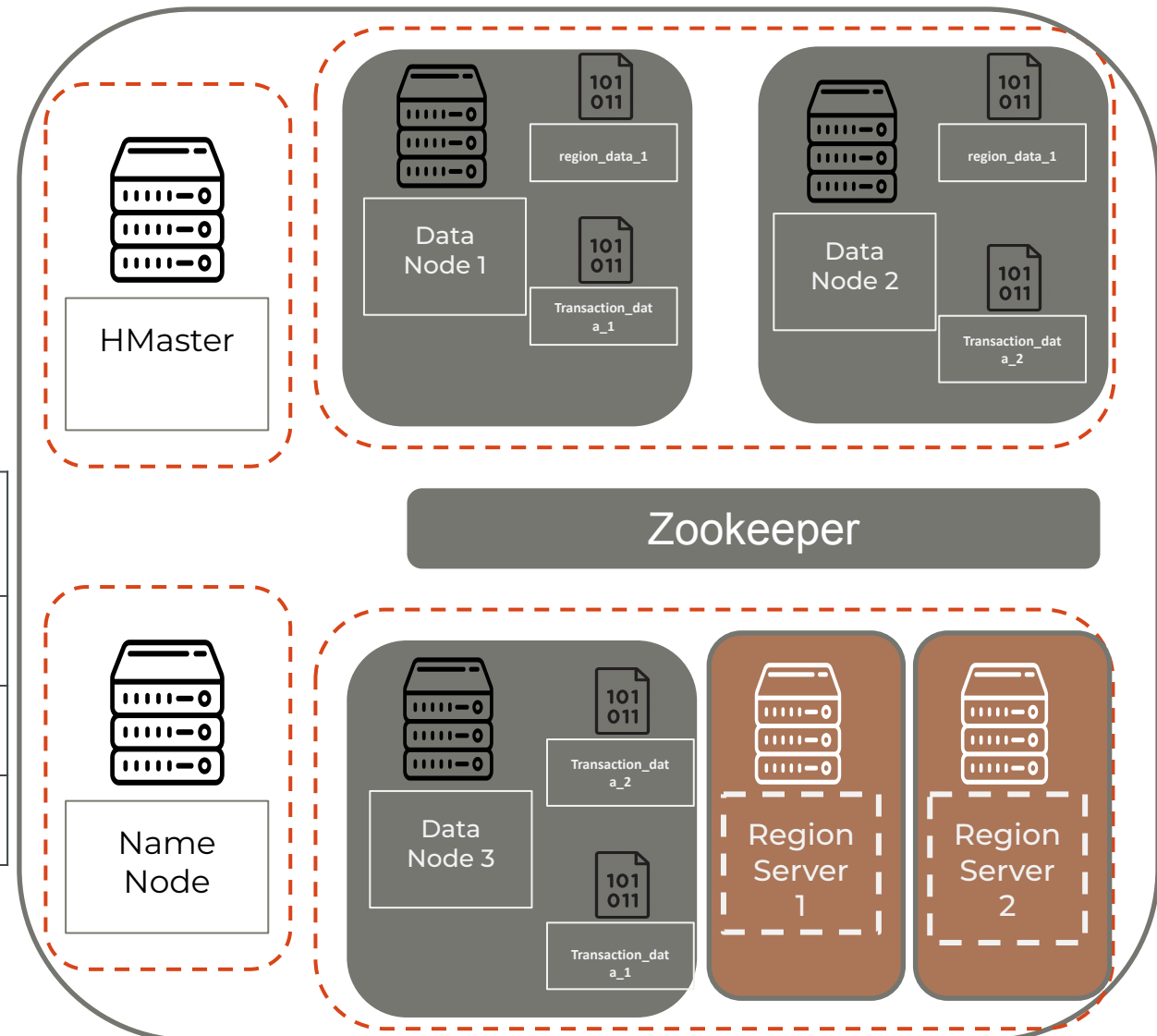
HBase Funcionamiento



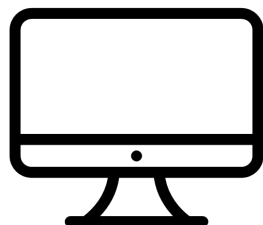
Client-1

```
hbase put
electric_company_cups_regiones
'0001_mun', 'Tarificacion_data:pico'
'11.34'
```

Row	Region_data	Tarificacion_data		timestamp
	valor_region	valle	pico	
0001_mun	Colmenar Viejo			1631952266000
0001_mun		10.34		1631952400000
0001_mun			11.34	1631952500000



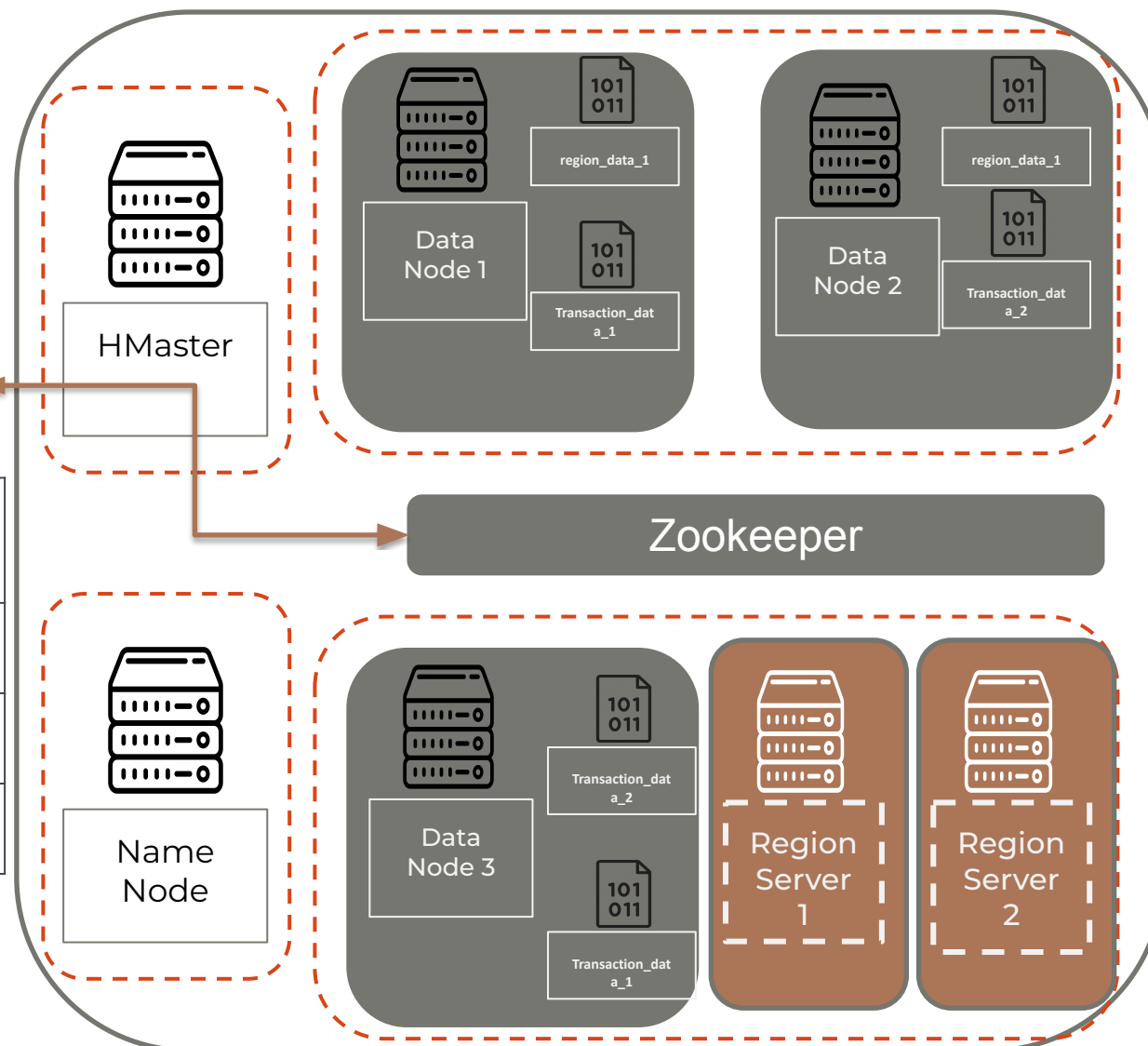
HBase Funcionamiento



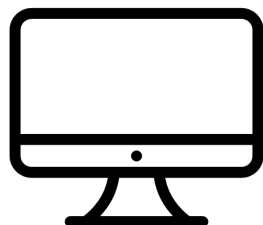
Client-1

```
hbase put
electric_company_cups_regiones
'0002_mun', region_data:valor_region
'Alcalá'
```

Row	Region_data	Tarificacion_data		timestamp
	valor_region	valle	pico	
0001_mun	Colmenar Viejo			1631952266000
0001_mun		10.34		1631952400000
0001_mun			11.34	1631952500000



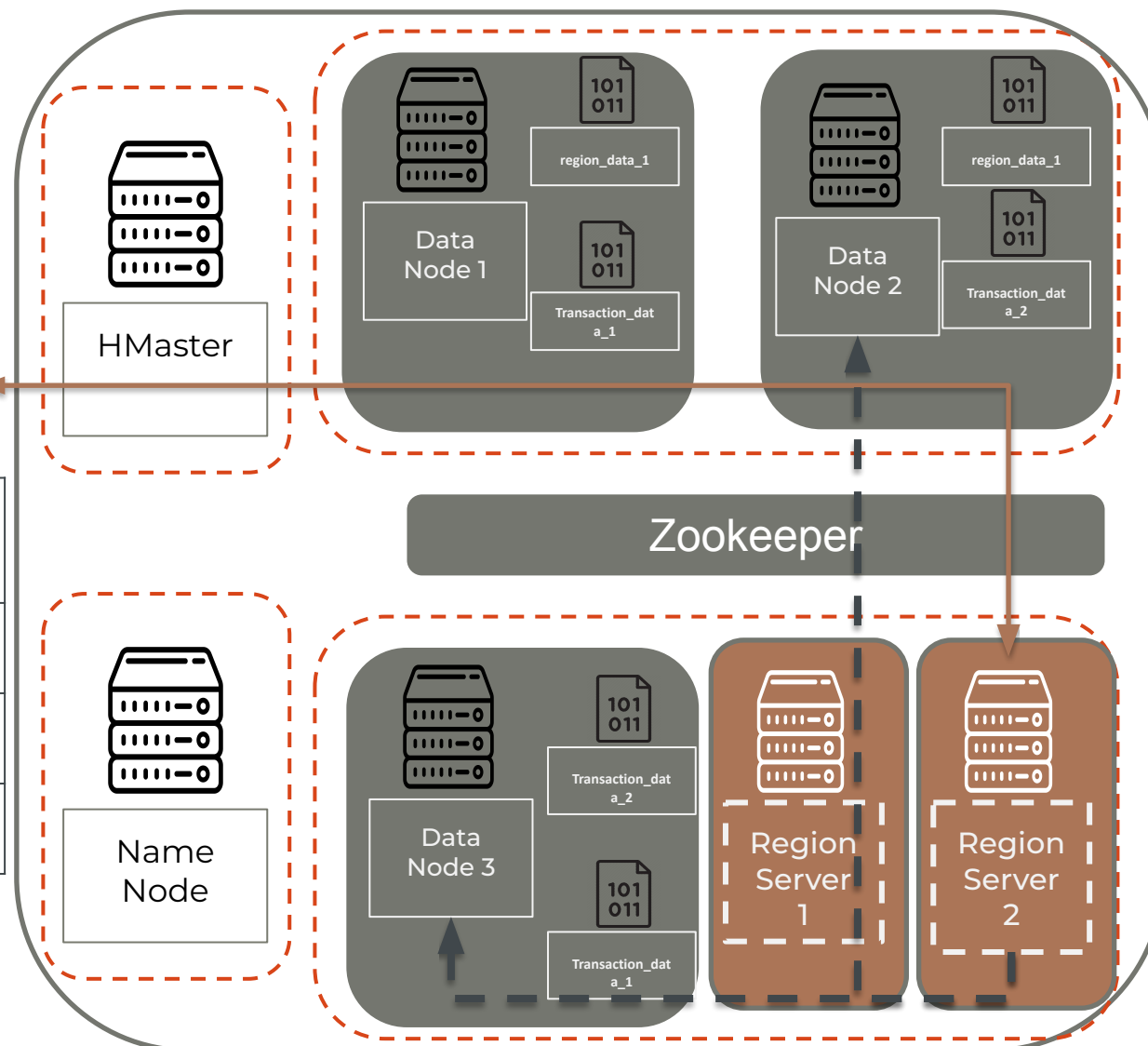
HBase Funcionamiento



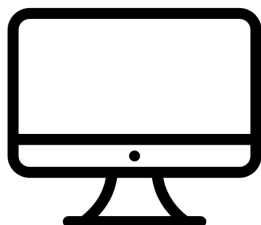
Client-1

```
hbase put
electric_company_cups_regiones
'0002_mun', region_data:valor_region
'Alcalá'
```

Row	Region_data	Tarificacion_data		timestamp
	valor_region	valle	pico	
0001_mun	Colmenar Viejo			1631952266000
0001_mun		10.34		1631952400000
0001_mun			11.34	1631952500000



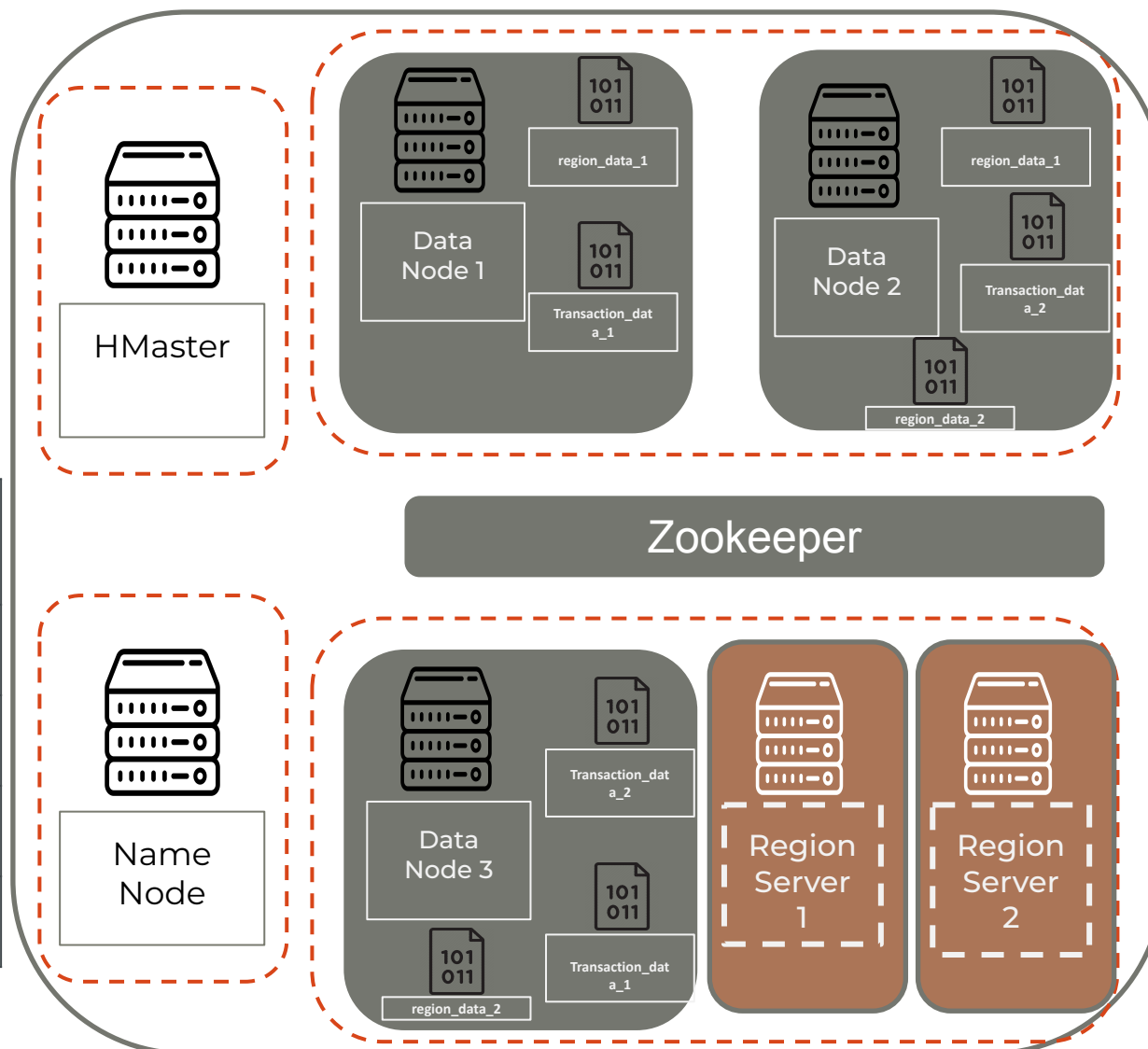
HBase Funcionamiento



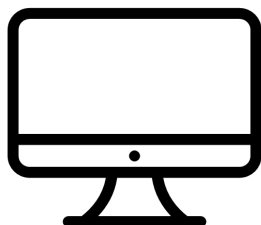
Client-1

```
hbase put
electric_company_cups_regiones
'0002_mun', region_data:valor_region
'Alcalá'
```

Row	Region_data	Tarificacion_data		timestamp
	valor_region	valle	pico	
0001_mun	Colmenar Viejo			1631952266000
0001_mun		10.34		1631952400000
0001_mun			11.34	1631952500000
0002_mun	Acalá			1631952600000



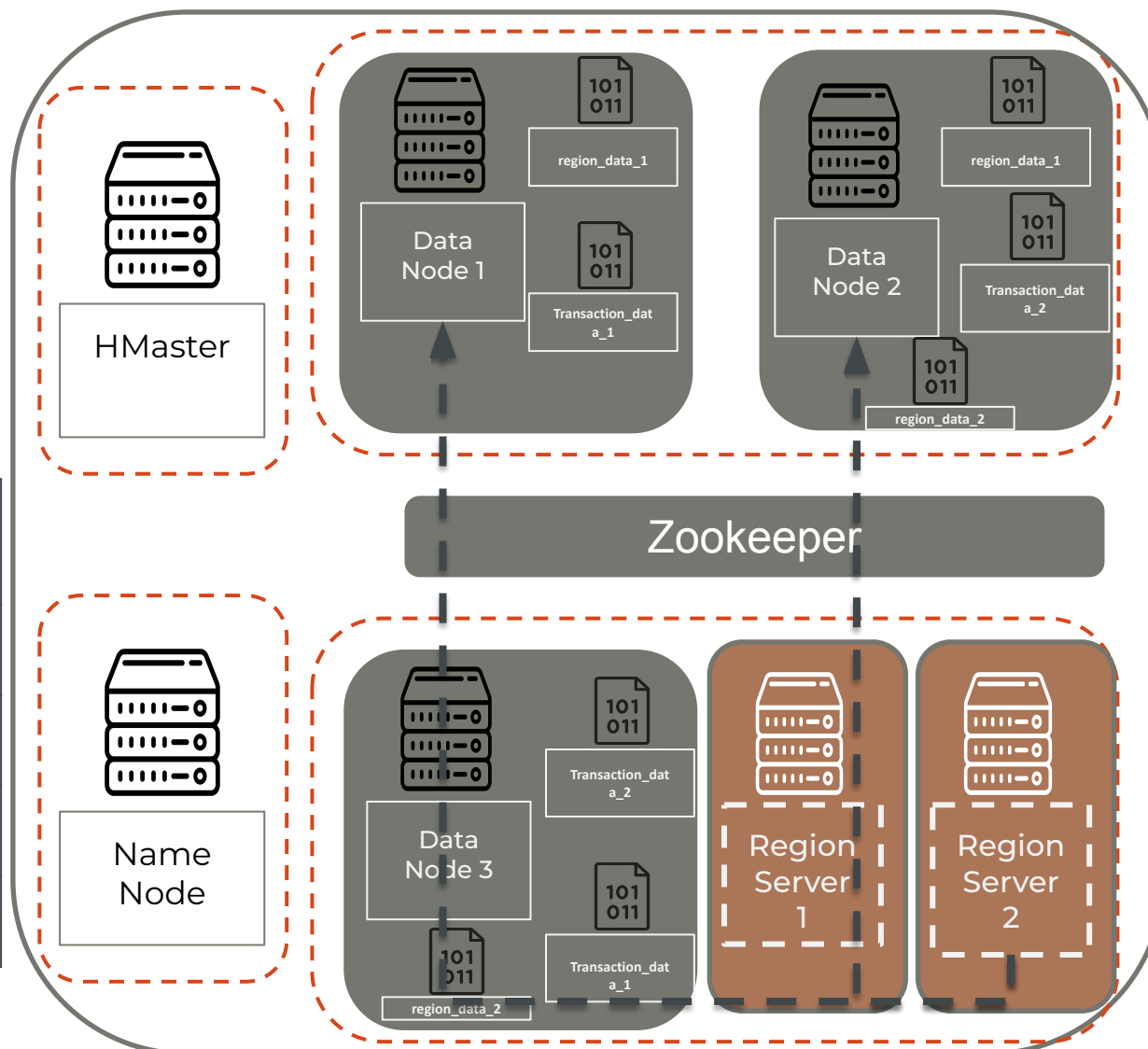
HBase Funcionamiento



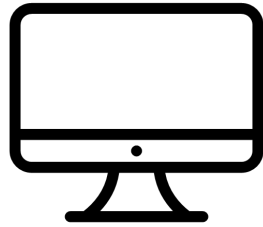
Client-1

```
hbase put
electric_company_cups_regiones
'0002_mun', region_data:valor_region
'Alcalá'
```

Row	Region_data	Tarificacion_data		timestamp
	valor_region	valle	pico	
0001_mun	Colmenar Viejo			1631952266000
0001_mun		10.34		1631952400000
0001_mun			11.34	1631952500000
0002_mun	Acalá			1631952600000

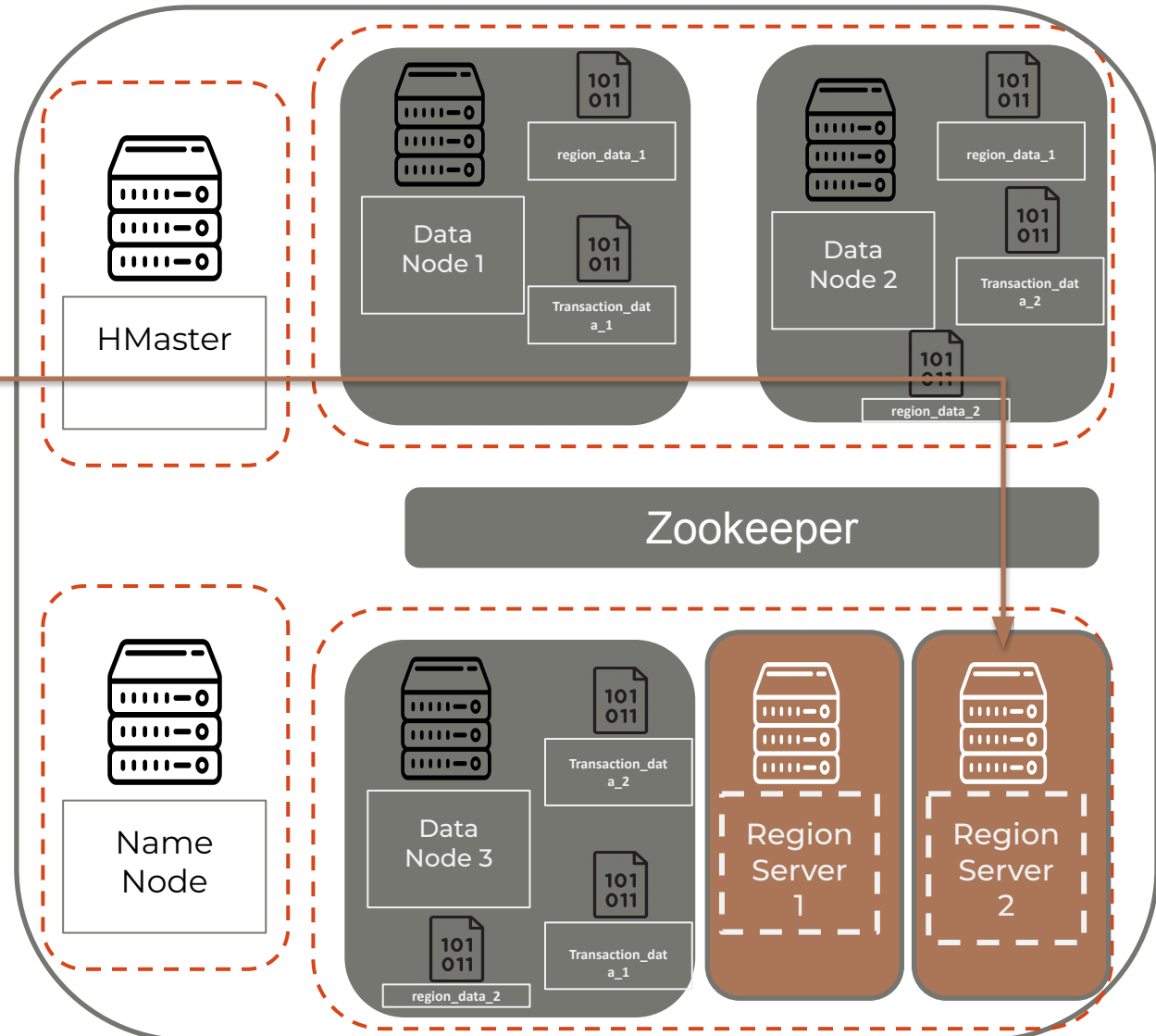


HBase Funcionamiento

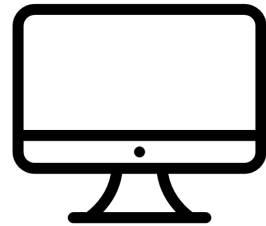


Client-1

```
hbase get  
electric_company_cups_regiones,  
'0002_mun'
```

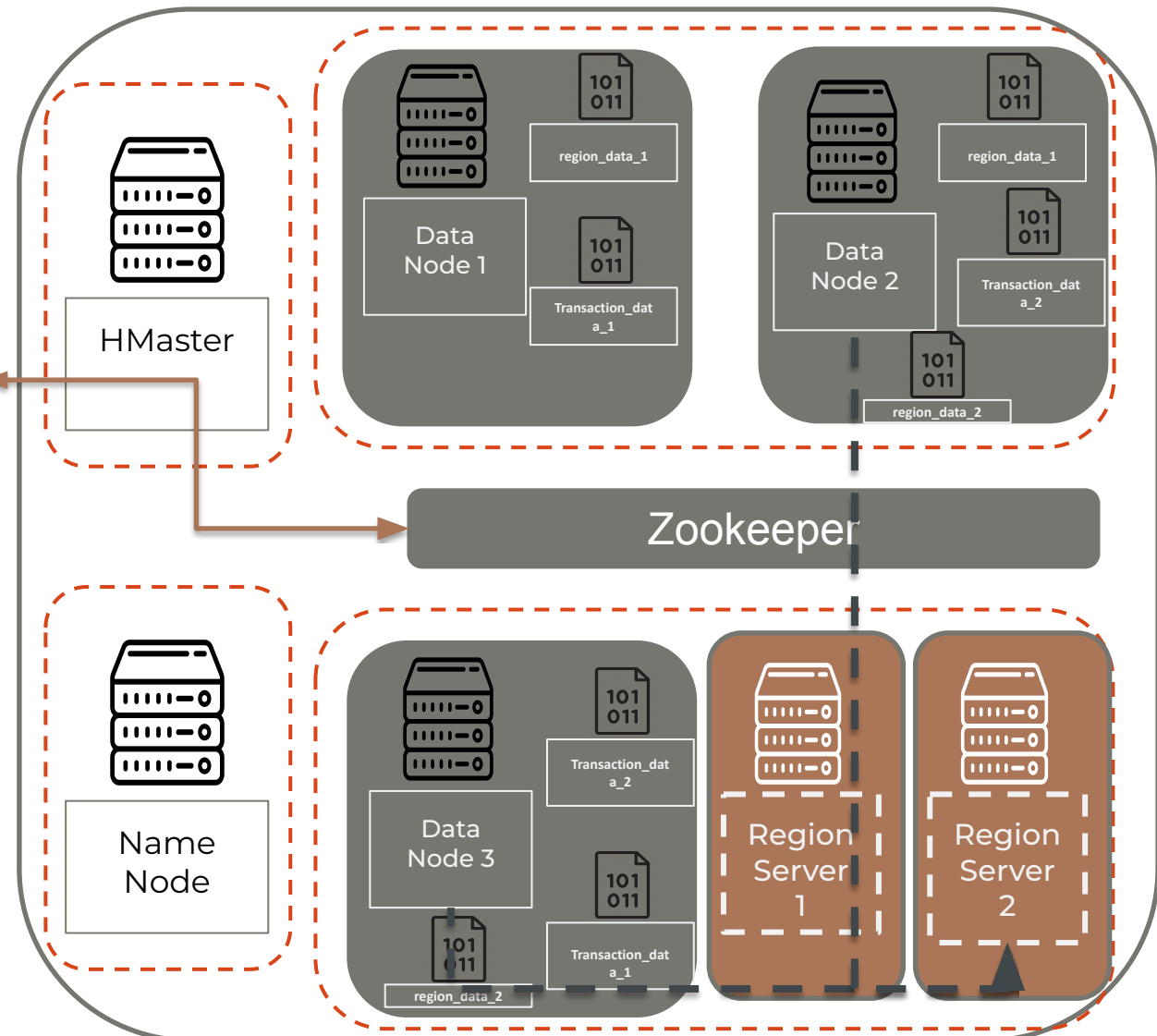


HBase Funcionamiento

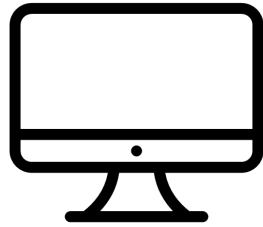


Client-1

```
hbase get
electric_company_cups_regiones,
'0002_mun'
```

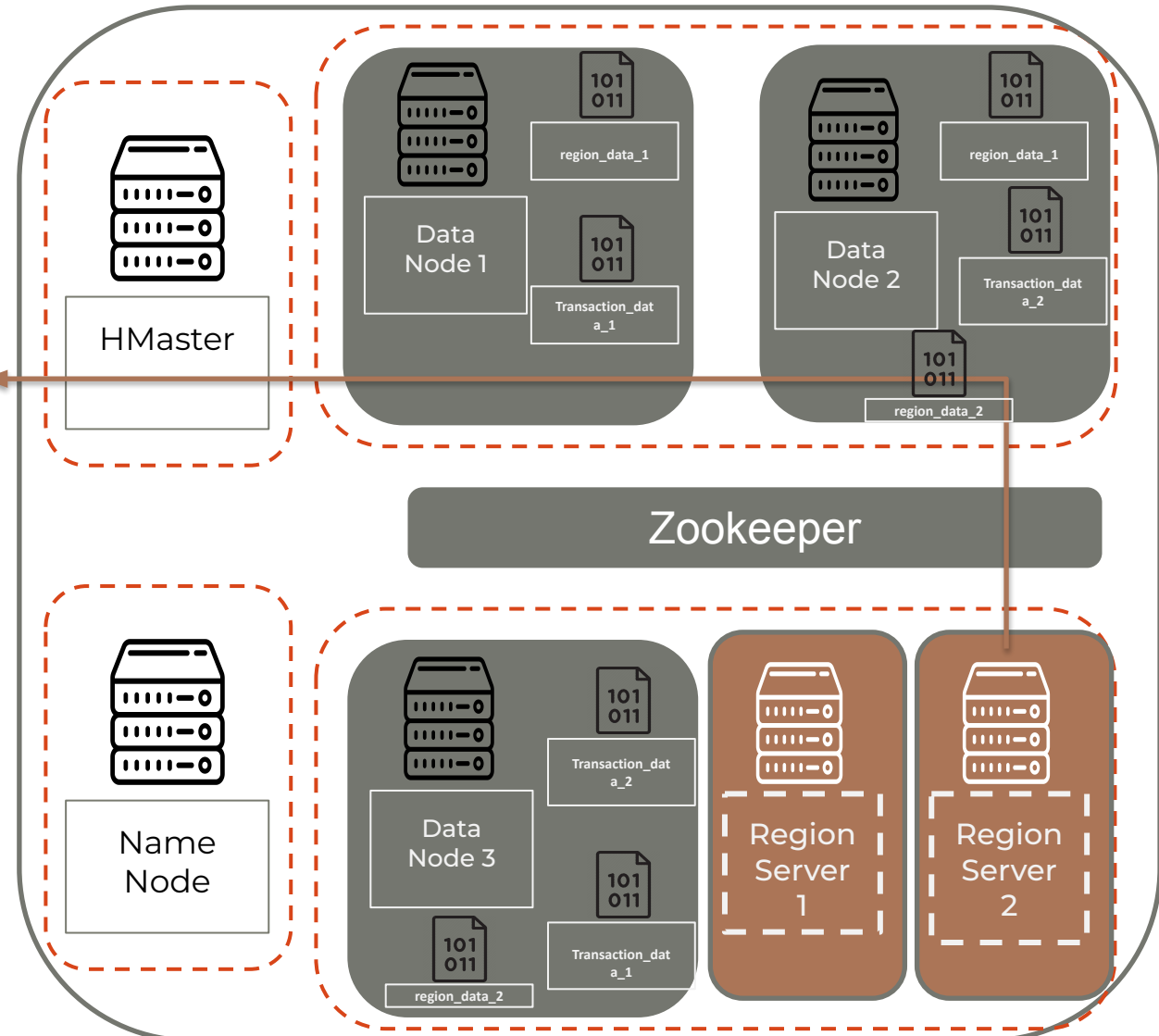


HBase Funcionamiento

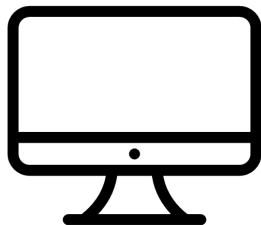


Client-1

```
hbase get  
electric_company_cups_regiones,  
'0002_mun'
```



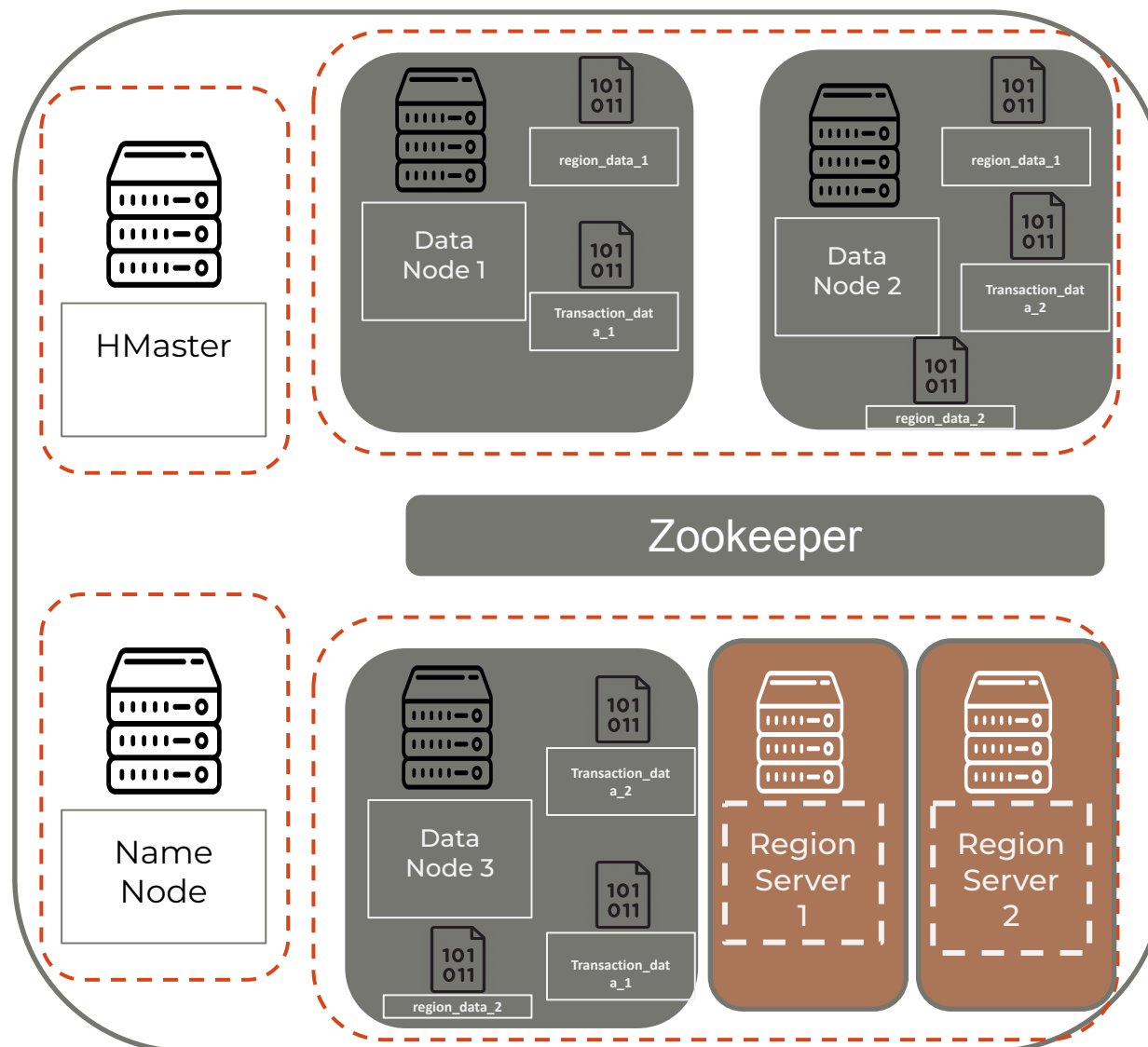
HBase Funcionamiento



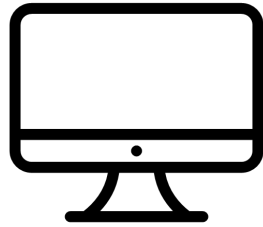
Client-1

```
hbase get
electric_company_cups_regiones,
'0002_mun'
```

COLUMN	CELL
region_data	valor_region
timestamp=1631952266000,	value='Colmenar Viejo'
Tarifificacion_data:pico	
timestamp=1631952500000,	value=11.34
Tarifificacion_data:valle	
timestamp=1631954000000,	value=100.34

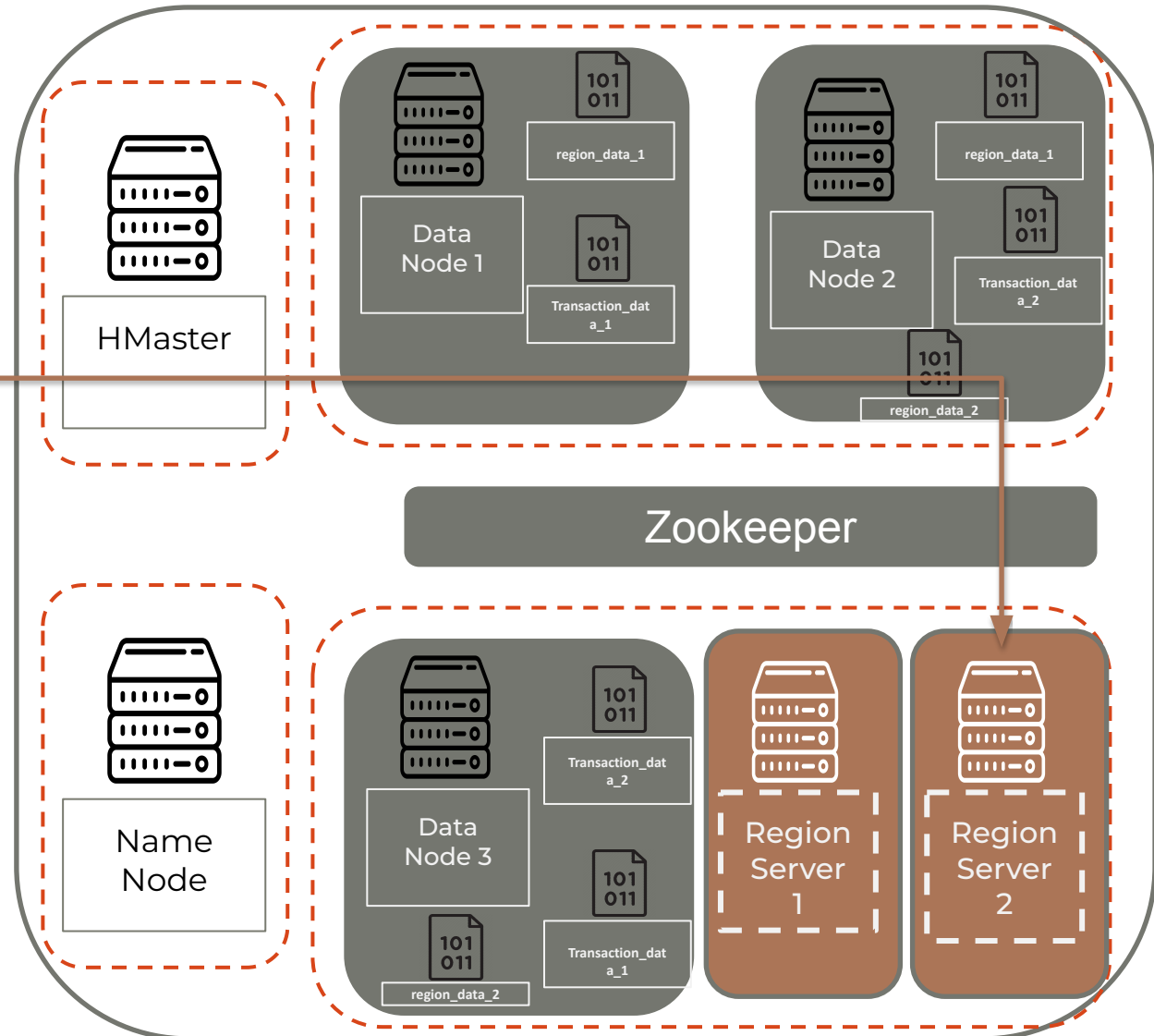


HBase Funcionamiento

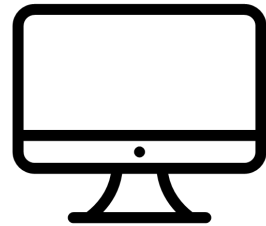


Client-1

```
hbase get
electric_company_cups_regiones,
'0002_mun'
```

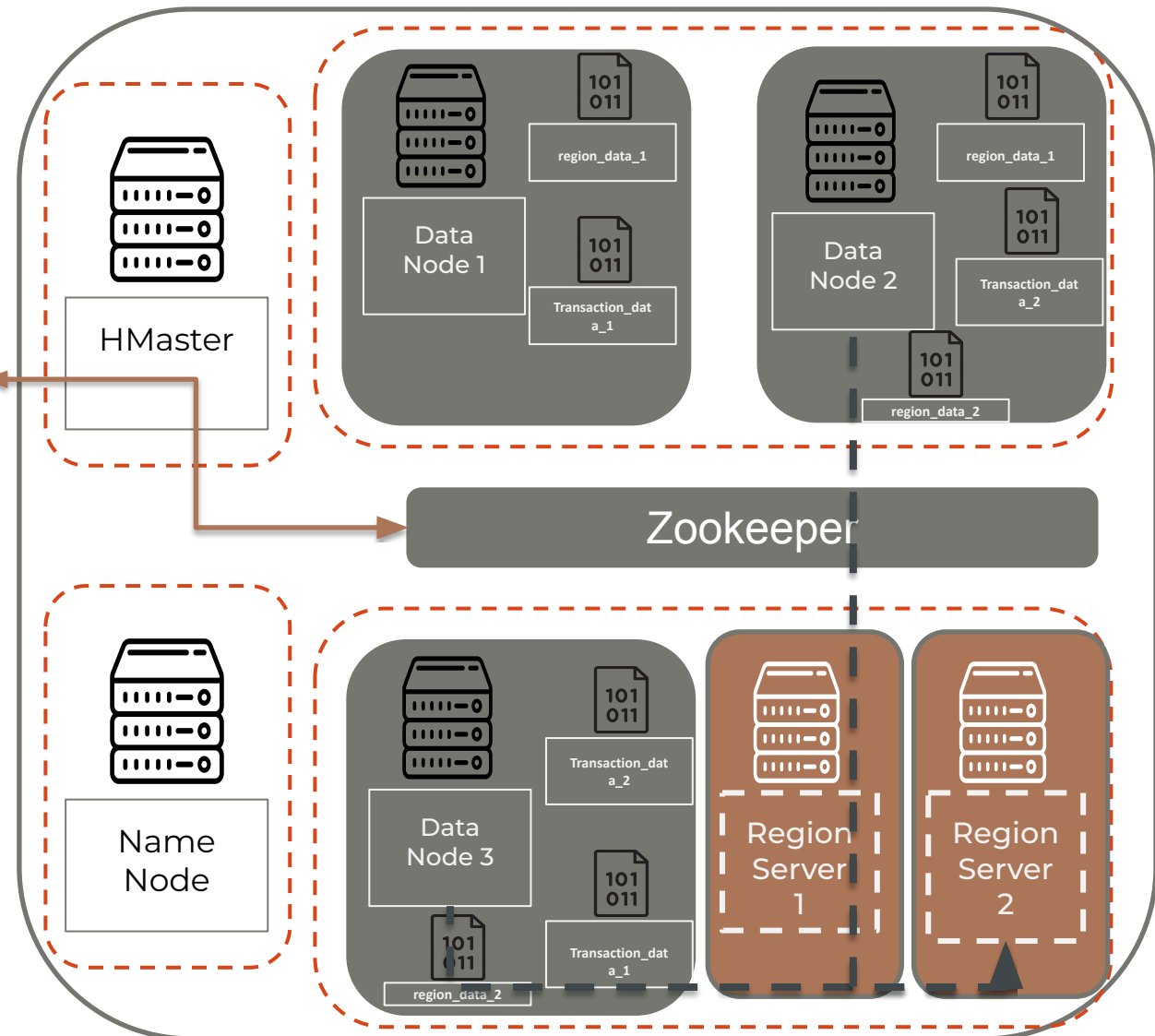


HBase Funcionamiento

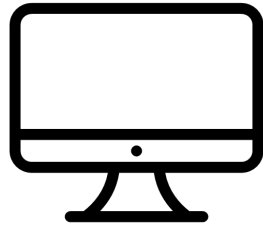


Client-1

```
hbase get
electric_company_cups_regiones,
'0002_mun'
```

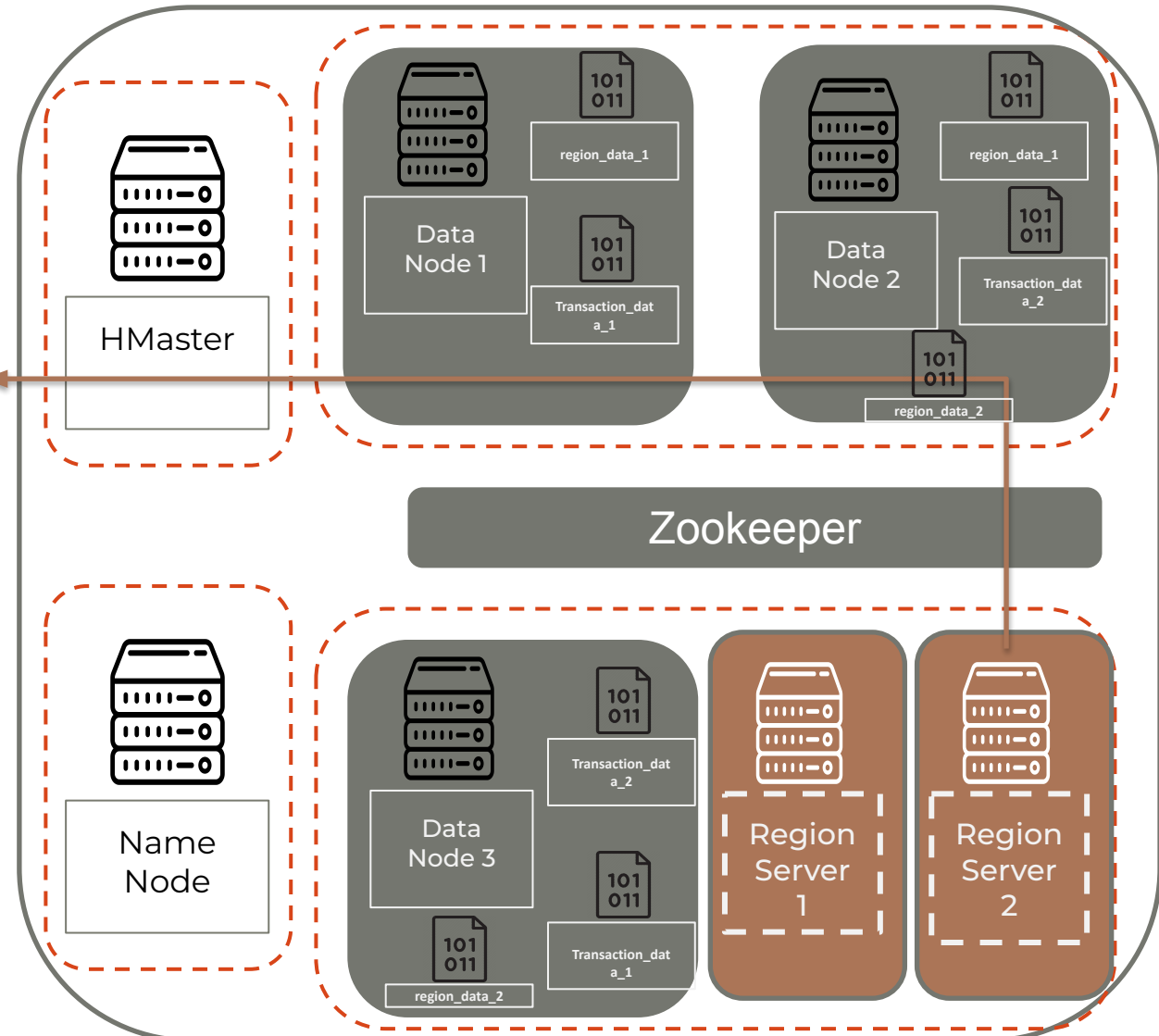


HBase Funcionamiento

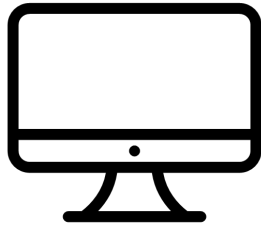


Client-1

```
hbase get  
electric_company_cups_regiones,  
'0002_mun'
```

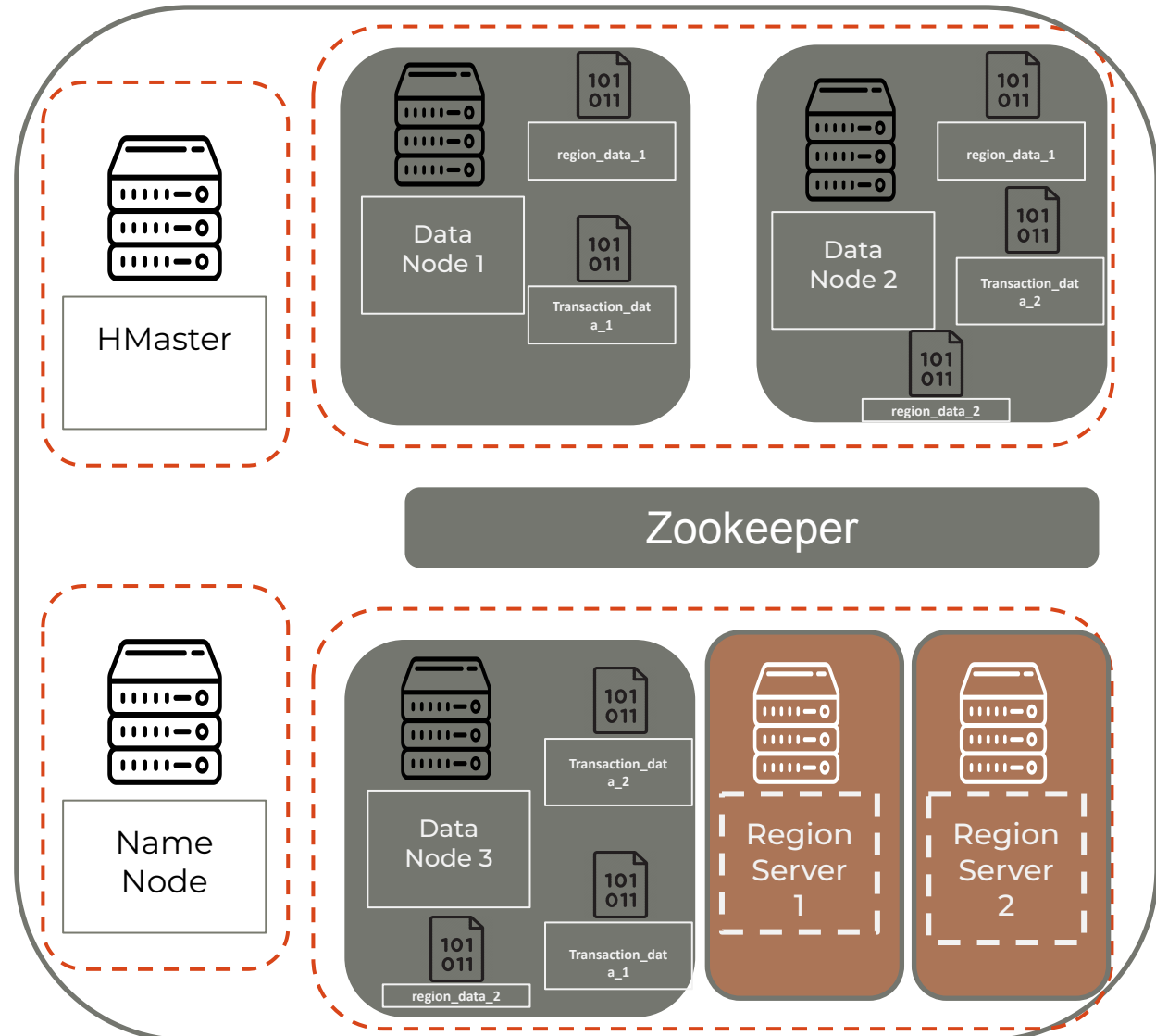


HBase Funcionamiento

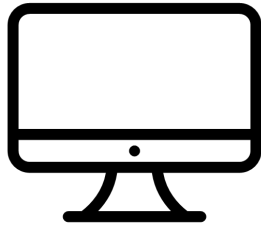


Client-1

```
hbase scan  
electric_company_cups_regiones
```

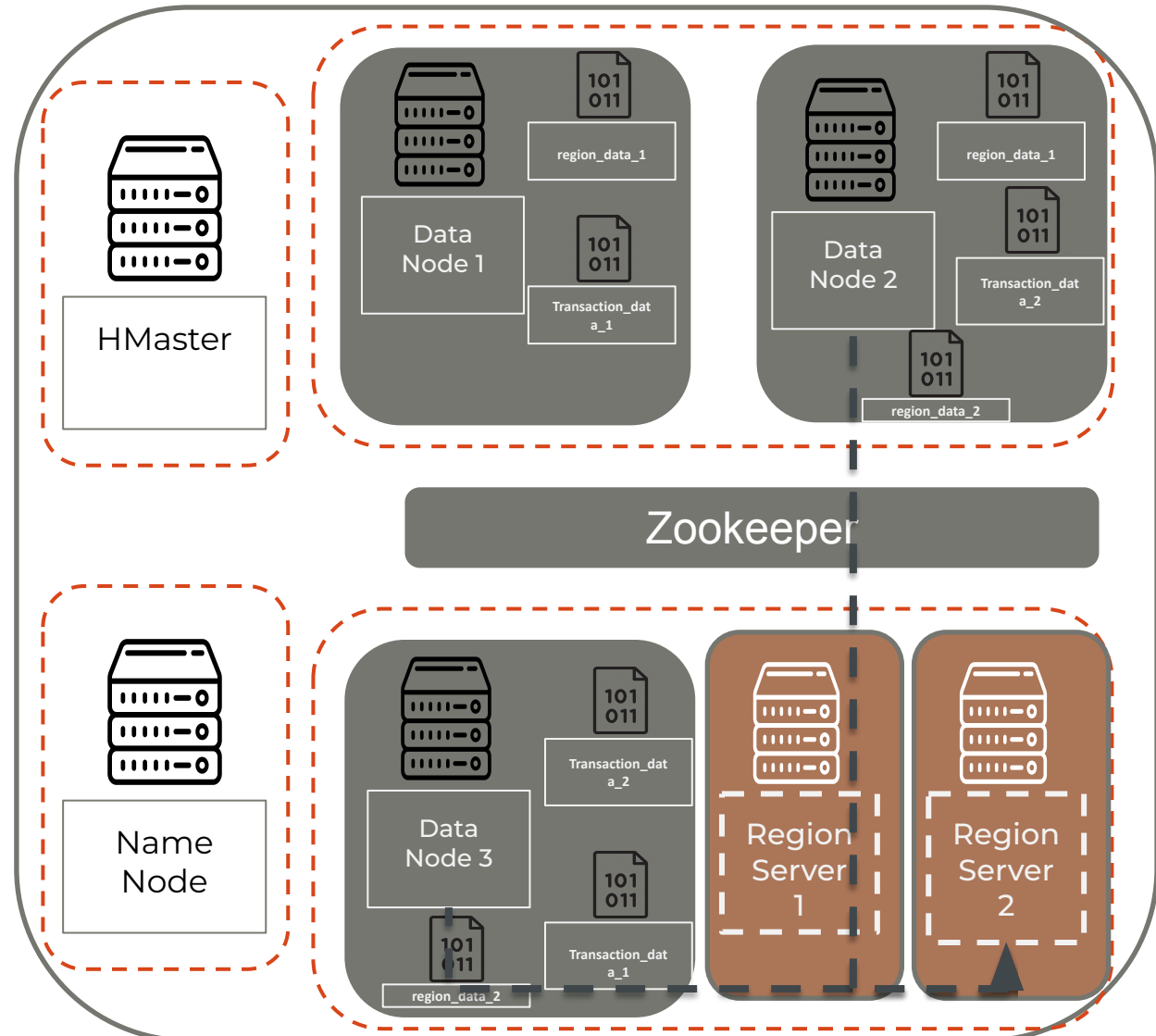


HBase Funcionamiento

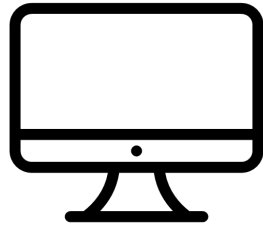


Client-1

```
hbase scan  
electric_company_cups_regiones
```

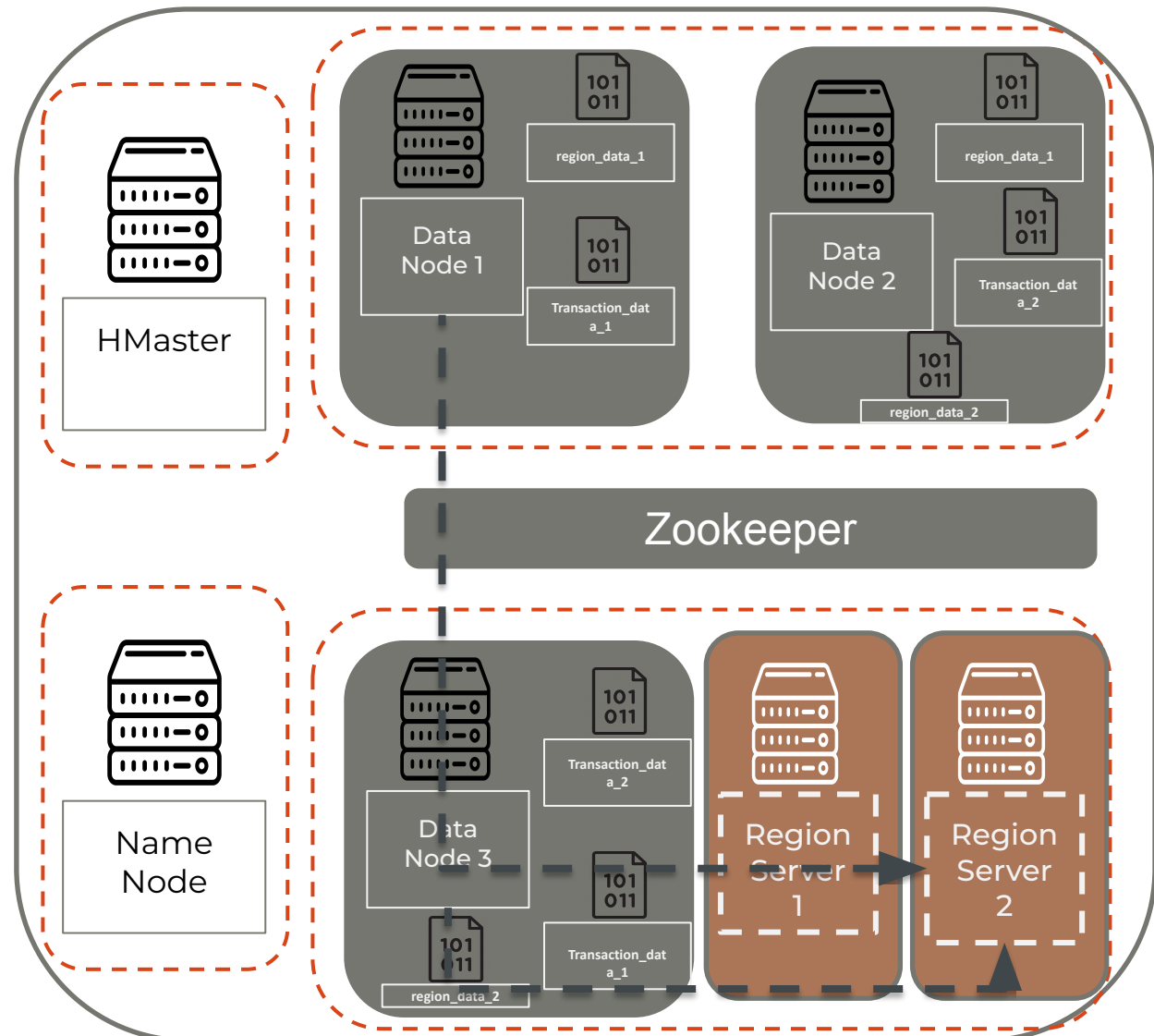


HBase Funcionamiento

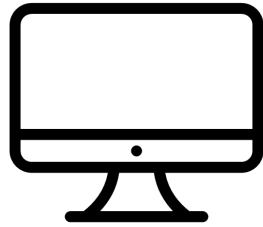


Client-1

```
hbase scan  
electric_company_cups_regiones
```

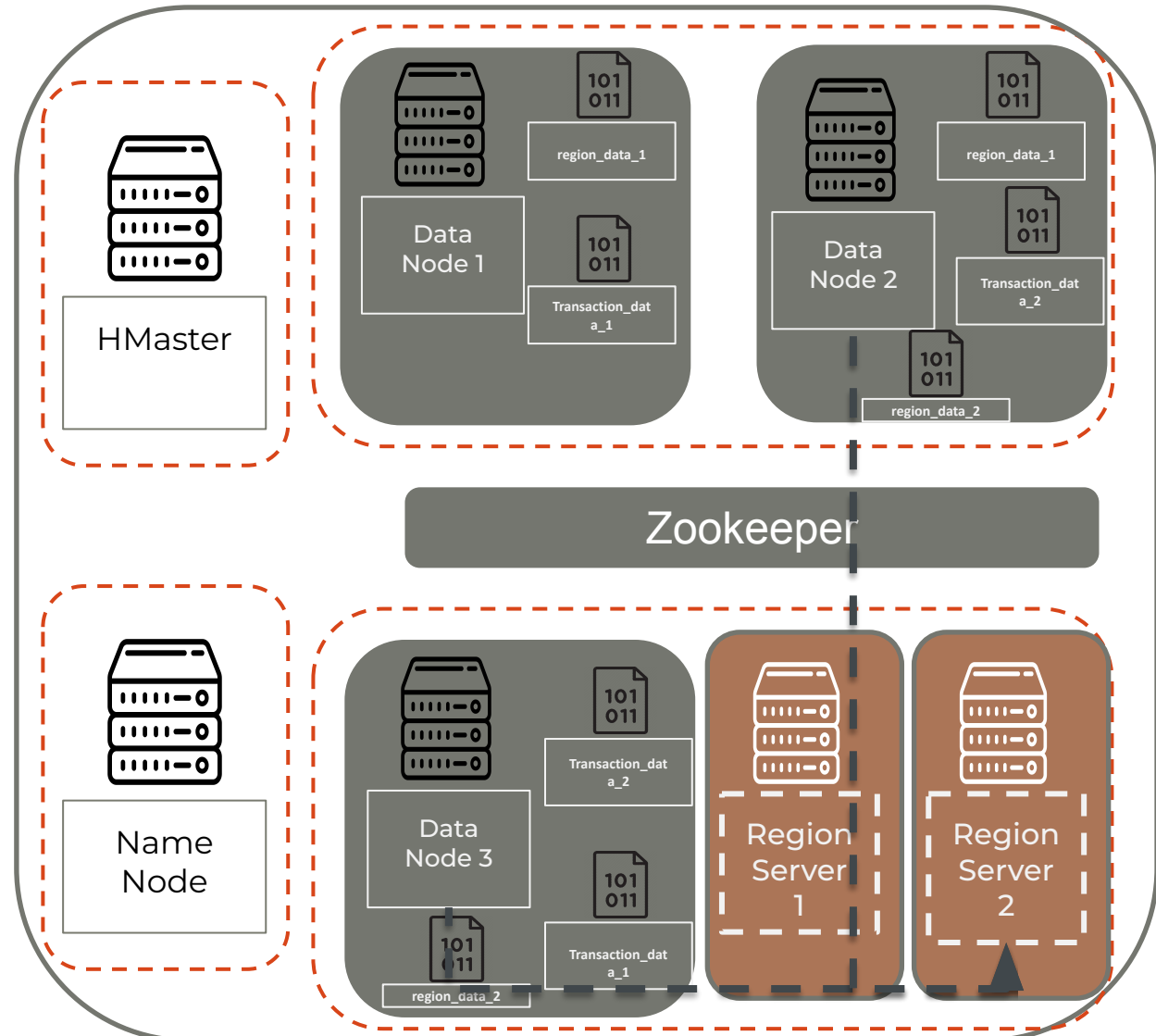


HBase Funcionamiento

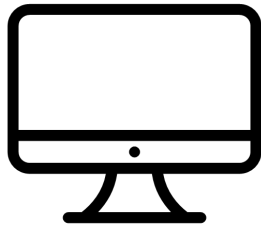


Client-1

```
hbase scan  
electric_company_cups_regiones
```



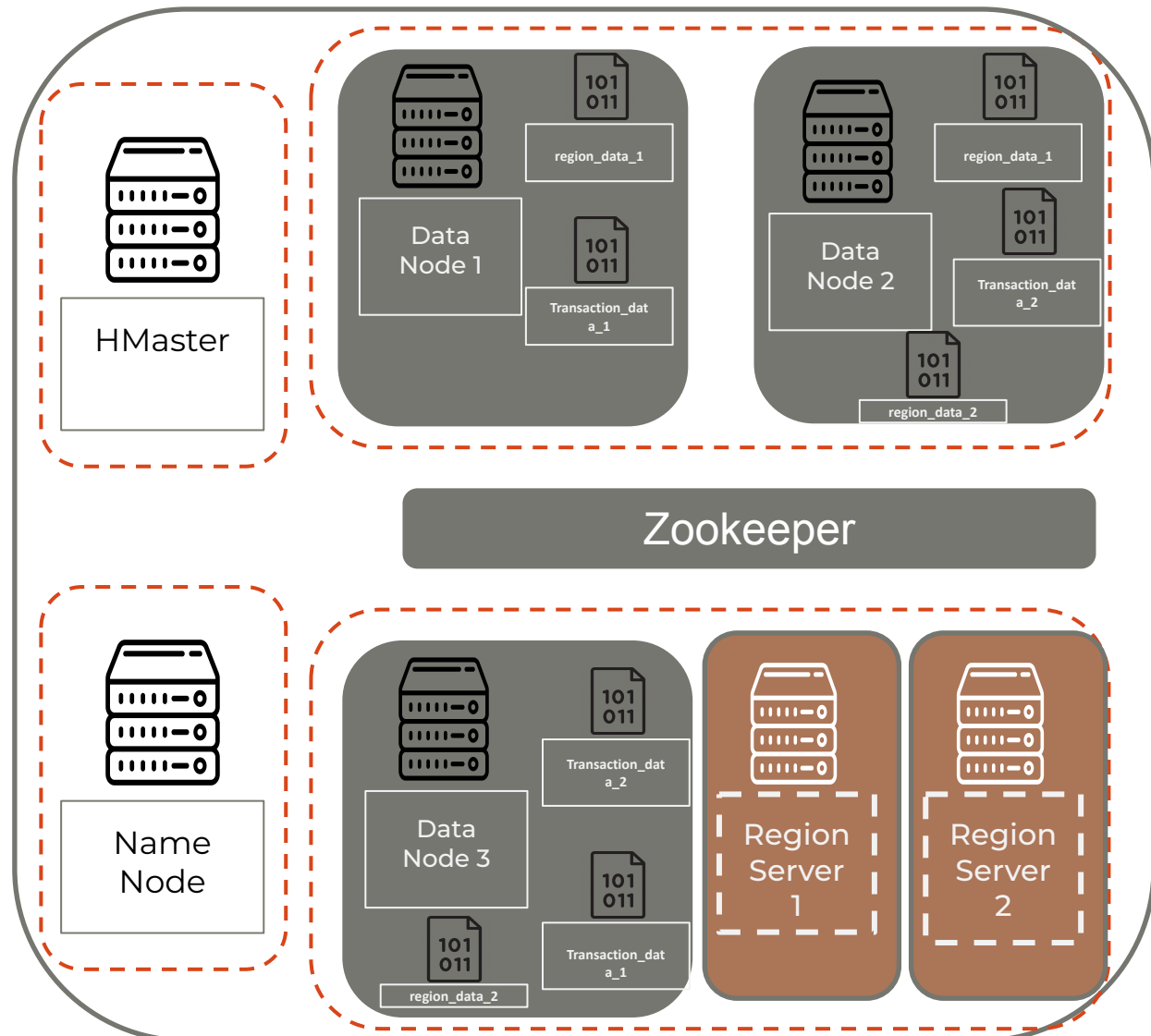
HBase Funcionamiento



Client-1

```
hbase scan  
electric_company_cups_regiones
```

```
ROW COLUMN+CELL  
  '0001-municipio'  
region_data : valor_region  
timestamp=1631952266000,  
value='Colmenar Viejo'  
  '0001-municipio'  
Tarificacion_data:pico  
timestamp=1631952500000, value=11.34  
  '0001-municipio'  
Tarificacion_data:valle  
timestamp=1631954000000, value=100.34  
  '0002-municipio'  
region_data : valor_region  
timestamp=1631952600000,  
value='Alcalá'
```



HBase - Introducción

- A raíz de las nuevas formas de almacenar se empiezan a pensar nuevas formas de consultar esa información
- El principal problema que pretende resolver es el acceso aleatorio a unos datos concretos
- Para ello usa HDFS como sistema de ficheros subyacente y una estrategia **Columnar**
- Con ello se consigue un acceso “rápido” a datos aleatorios en colecciones masivas de información
- Posteriormente nacen tecnologías de almacenamiento basadas en su misma filosofía supliendo algunas de sus carencias (Apache Cassandra)



Afi Escuela

© 2022 Afi Escuela. Todos los derechos reservados.