

# Procesamiento de Datos Masivos

## Tarea 1

### 1. Esquema de datos

#### Parlamentarios

```
1 CREATE OR REPLACE TABLE proceso-de-datos-454919.tarea1.parlamentarios AS
2 SELECT
3     p.PARLAMENTARIO_ID AS id_parlamentario,
4     p.NOMBRE_COMPLETO AS nombre_completo,
5     pt.id_partido AS partido_id
6 FROM proceso-de-datos-454919.tarea1.parlamentarios_raw p
7 LEFT JOIN proceso -de-datos-454919.tarea1.partidos pt
8 ON p.PARTIDO_POLITICO = pt.nombre_partido;
```

#### Partidos

```
1 CREATE OR REPLACE TABLE proceso -de-datos-454919.tarea1.partidos AS
2 SELECT
3     DENSE_RANK() OVER (ORDER BY PARTIDO_POLITICO) AS id_partido,
4     PARTIDO_POLITICO AS nombre_partido
5 FROM (
6     SELECT DISTINCT PARTIDO_POLITICO
7     FROM proceso -de-datos-454919.tarea1.parlamentarios_raw
8     WHERE PARTIDO_POLITICO IS NOT NULL
9 );
```

#### Keywords

```
1 CREATE OR REPLACE TABLE proceso-de-datos-454919.tarea1.keywords AS
2 SELECT
3     ROW_NUMBER() OVER () AS id_keyword,
4     palabra
5 FROM (
6     SELECT DISTINCT palabra
7     FROM proceso -de-datos-454919.tarea1.dataframe,
8         UNNEST(
9             SPLIT(
10                 REPLACE(REPLACE(REPLACE(keywords, "[", ""), "]", ""), "'", ""),
11                 ",", ""
12             )
13         ) AS palabra
14 );
```

## Intervenciones

```
1 CREATE OR REPLACE TABLE proceso -de-datos-454919.tarea1.intervenciones
2 AS
3 SELECT
4     DISTINCT intervention_id AS id,
5     p.id AS parlamentario_id,
6     intervention_date AS fecha
7 FROM proceso -de-datos-454919.tarea1.parlamentarios_raw r
8 JOIN proceso -de-datos-454919.tarea1.parlamentarios p
   ON r.NOMBRE_COMPLETO= p.nombre_completo;
```

## Intervenciones-keywords

```
1 CREATE OR REPLACE TABLE
2 proceso -de-datos-454919.tarea1.intervenciones_keywords AS
3 SELECT
4     d.ID_PARTICIPACION AS intervencion_id,
5     k.id_keyword AS keyword_id
6 FROM proceso -de-datos-454919.tarea1.dataframe d,
7     UNNEST(
8         SPLIT(
9             REPLACE(REPLACE(REPLACE(d.keywords, "[", ""), "]", ""), "'", ""),
10             ", "
11         )
12     ) AS palabra
13 JOIN proceso -de-datos-454919.tarea1.keywords k
14     ON k.palabra = palabra;
```

## 2. Modelación

En la figura 1 se puede ver el diagrama de modelación.

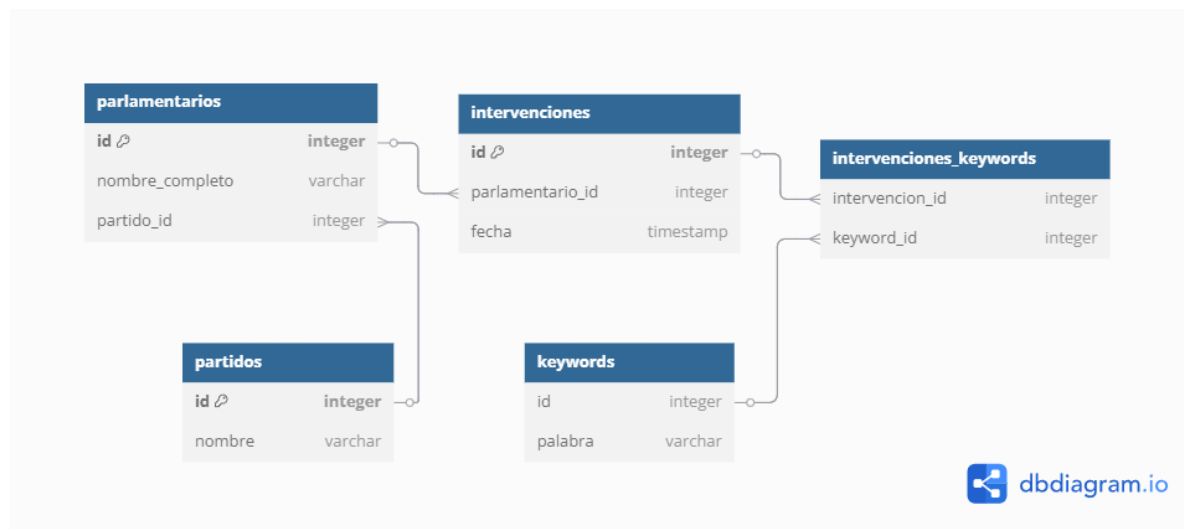


Figura 1: Diagrama de modelación

## Justificación

Se compone de cinco tablas principales: partidos, parlamentarios, intervenciones, keywords e intervenciones\_keywords. La estructura permite registrar qué parlamentario (y por ende qué partido) realizó una intervención en una fecha determinada, y vincular esa intervención con una o más palabras clave o temáticas a través de una tabla intermedia. Esto facilita, por ejemplo, consultas como identificar las 5 temáticas más tratadas en cada mes, ya que se puede agrupar por mes utilizando la columna de fecha en intervenciones, y luego contar la frecuencia de cada keyword mediante la relación con intervenciones\_keywords. También permite calcular la media móvil de intervenciones por partido político usando la fecha de cada intervención y asociándola al partido a través del parlamentario. Además, gracias a esta estructura relacional, es sencillo obtener, para cada trimestre, el tema más tratado por cada partido, o pararse en un mes específico y consultar el top 3 de temáticas tratadas por partido. En resumen, se trata de un modelo normalizado, flexible y eficiente para consultas temporales y temáticas, ideal para análisis político y de discurso parlamentario