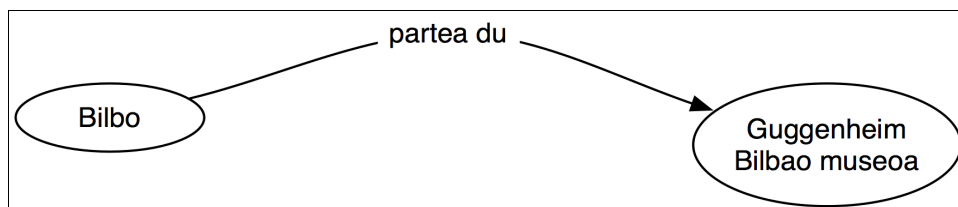


# DATUEN AMARAUNA

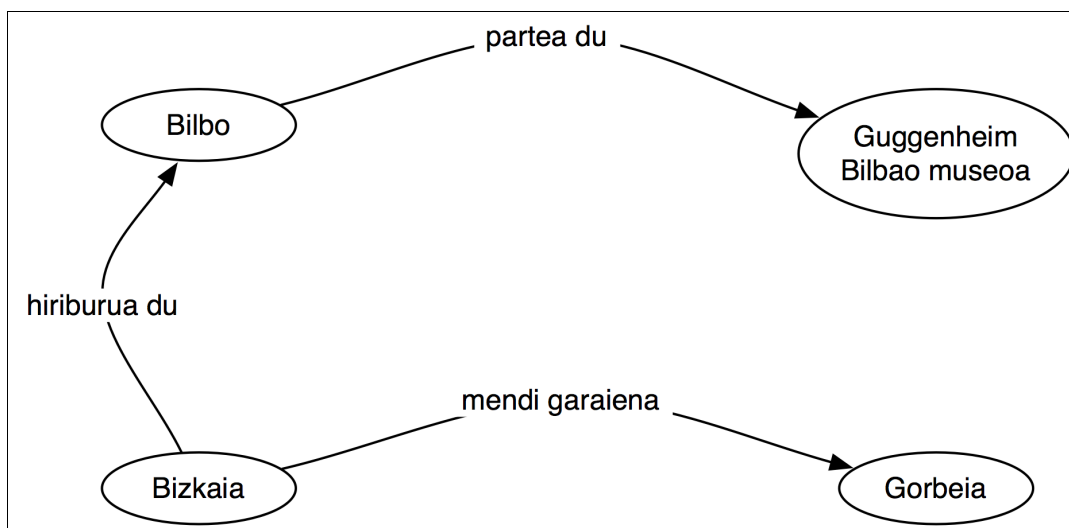
## 1.- INFORMAZIOA INTERNETEN ZUZENEAN ARGITARATU

Web sarean gabiltzanean, gure xedea informazioa aurkitzea da. Hala ere, informazioa eskuratu ordez, informazioa gordetzen duten web orrialdeetan hura bilatzen dugu, denbora galduz. Izan ere, web orrialdeetako testuak irakurri eta beharrezko informazioa erauzten dugu, informazioa zuzenean jo beharrean. Datu estekatuen teknologiak (*Linked Data*) konponbidea eskaintzen digu [1]: datu estekatuak interneten zuzenean argitara ditzakegu, web orrialdeak eraiki beharrik gabe, beste erabiltzaile batzuk datuok nahi bezala erabil ditzaten.

Datu estekatuen funtsa RDF (*Resource Description Framework*) lengoia da. RDFk informazioa egitura jakin batean adierazten du: RDF triplea (1 irudia). Hiru elementuk osatzen dute RDF triplea: subjektua, predikatua eta objektua. Predikatuak subjektua eta objektua lotzen ditu. Adibidez, *Bilbok parte du Guggenheim Bilbao museoa* triplean *Bilbo* subjektua *parte du* predikatuaren bidez *Guggenheim Bilbao museoa* objektuarekin erlazionatzen da. Beraz, predikatuak bi entitateen arteko erlazioak dira: entitate batek (Bilbok), propietate jakin batean (*Zein parte duen*) beste entitate bat du (*Guggenheim Bilbao museoa*). Horrelako tripleak lotuz informazio konplexua adierazten duten RDF sareak eratzen dira (2 irudia).



1 irudia: RDF triple baten adibidea (Subjektua: *Bilbo*; Predikatua: *parte du*; objektua: *Guggenheim Bilbao museoa*). Triple egituran edozein motako informazioa adieraz daiteke.



2 irudia: RDF tripleak lotuz eratzen den RDF sarea (Triple baten subjektua beste triple baten objektua izan daiteke: *Bilbo*, adibidez). RDF sarea informazio konplexua adierazteko erabiltzen da.



## 2.- DATUEN AMARAUNAREN ERABILERA

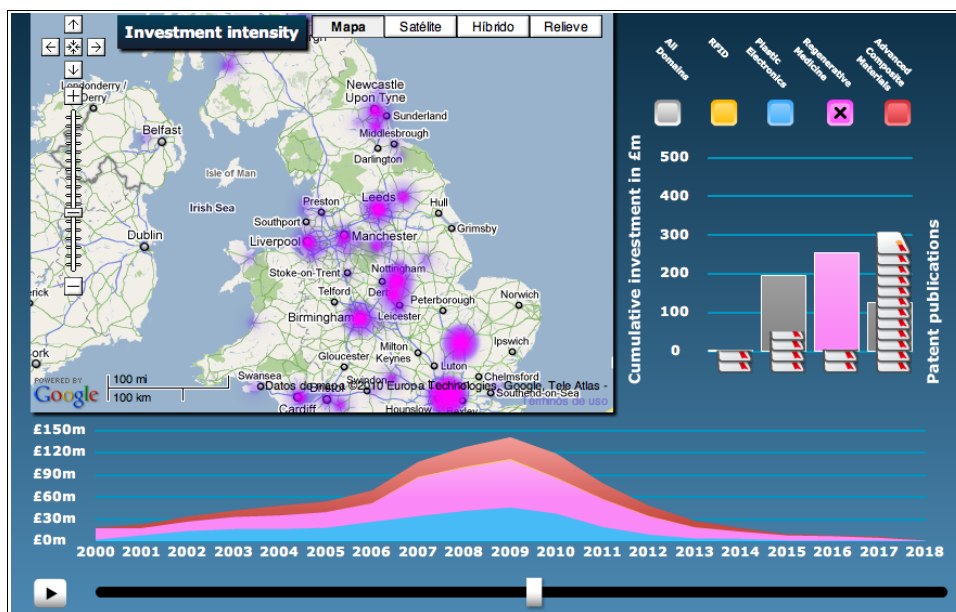
Datu estekatuetan dagoen informazioa erabiltzeko software programen ekosistema aberatsa sortu da, batez ere bi arrazoiengatik. Batetik, ohizko web sarearekin alderatuz datu estekatu bidez informazio zehatzagoa arinago lor dezakegu. Bestetik, informazio horrekin lan egingo duen softwarea aise garatzea dago, informazioa RDF bidez egituratua egongo denez, software garatzaileak aldeztu aurretik jakingo duela den informazio horren egitura (Hainbat RDF triplez osatutako RDF sarea).

Datuen amaraunari datu berriak erantsi erraza da: Triple egitura jarraitzea besterik ez dugu, eta gure tripleak predikatu bidez beste norbaiten tripleekin estekatu. Ondorioz, datuen amarauna egunero hedatzen da. Erabiltzaileek datu berri horiekiko elkarrekintza interesgarriak izan dezaten software programak garatzen dira egunero ere. Datu gehiago ez ezik, datu mota berriek software programa ere interesgarriagoak garatzea eskatuko dute etorkizunean.

### 2.1.- GOBERNUEN DATU IREKIAK

Hiritarrok gobernuak gordetzen dituen datuekin lan egiteko eskubidea dugu. Hala ere, orain arte datu horiek datu-base itxietan gorde izan dira, hiritarren eskutik at. Gaur egun, RDFk datuok argitaratzeko teknologia paregabea eskaintzen du: datuak beste edozein daturekin uztartzeko moduan argitara daitezke, eta hiritarrok, datuak aztertu ez ezik, software programak eraiki ditzakegu datuok lan egiteko. Ondorioz, gobernuak gero eta datu publiko gehiago argitaratzen dituzte datu estekatuaren ereduari jarraiki, hiritarrekiko gardentasuna bermatzeko.

Erresuma Batuko gobernuak<sup>1</sup> aitzindaria izan da datu estekatuak eskaintzeko prozesuan, eta ekimen horren inguruan software programa interesgarriak sortu dira. Adibidez, Talis enpresak ikerketa guneen mapa elkarreragilea garatu du<sup>2</sup>: mapa horretan Erresuma Batuko hiritarrek goi-mailako ikerketa guneak aurki ditzakete, lortutako diru eta patenteen arabera, kronologikoki ordenatuak (4 irudia).



4 irudia: Talis enpresak ikerketa gune garrantzitsuenak aurkitzeko sortutako mapa. Emaizak ikerketa arloaren arabera sailkatzeko aukera dago. [<http://bis.clients.talis.com/>]

1 <http://data.gov.uk>

2 <http://bis.clients.talis.com/>

## 2.2.- SENDAGAIK AURKITZEKO BIDE BERRIAK

Biologian informazio asko argitaratzen da interneten zehar barreiaturiko baliabideetan. Horrek esan nahi du ikerlari batek informazioa jaso eta bateratu behar duela galdera konkretuak ebazteko. Datu estekatuekin, aldiz, ikerlariak behar duen informazioa sistema bakarretik jaso dezake. Horrelako sistemetako bat Neurocommons da [2]: RDF bidez Alzheimer gaixotasunari buruzko informazio bateratua eskaintzen du. Ikerlariak Neurocommons sistema farmako berrientzako jomuga posibleak aurkitzeko erabiltzen dute. Adibidez, CA1 piramide neuronak oso kaltetuak daude Alzheimer gaixotasuna pairatzen duten pertsonetan, eta ikerlariak badakite neurona horietan seinale-transdukzioa garrantzitsua dela. Beraz, seinale-transdukzioa burutzen duten eta piramide neuronetan dauden proteinak jakinez gero, proteina horiek farmako berrientzako jomuga posibleak izango dira.

Galdera hori Googlen eginez gero (*Signal transduction in pyramidal neurons*) esanahi gabeko erantzunez jositako zerrenda handiegia lortuko dute ikerlariak. Aitzitik, Neurocommons sistema erabiliz (5 irudia), osagaiantzako jopuntu gutxi eta baliotsuak lortzen dituzte, Alzheimer gaixotasuna sendatzeko balizko tratamendu berriak, alegia.

```
1. prefix go: <http://purl.org/obo/owl/GO#>
2. prefix rdfs: <http://www.w3.org/2000/01/rdf-schema#>
3. prefix owl: <http://www.w3.org/2002/07/owl#>
4. prefix mesh: <http://purl.org/commons/record/mesh/>
5. prefix sc: <http://purl.org/science/owl/sciencecommons/>
6. prefix ro: <http://www.obofoundry.org/ro/ro.owl#>
7. SELECT ?gene_name ?process_name
8. WHERE
9. { ?pubmed_record ?related mesh:D017966 .
10.   ?article sc:identified_by_pmid ?pubmed_record.
11.   ?gene_record sc:describes_gene_or_gene_product_mentioned_by ?article.
12.   ?protein rdfs:subClassOf ?restriction.
13.   ?restriction owl:onProperty ro:has_function.
14.   ?restriction owl:someValuesFrom ?restriction2.
15.   ?restriction2 owl:onProperty ro:realized_as.
16.   ?restriction2 owl:someValuesFrom ?process.
17.   { { ?process ro:part_of go:GO 0007165 }
18.     UNION
19.     { ?process rdfs:subClassOf go:GO 0007165 }}
20.   ?protein rdfs:subClassOf ?parent.
21.   ?parent owl:equivalentClass ?restriction3.
22.   ?restriction3 owl:onProperty sc:is_protein_gene_product_of_dna_described_by.
23.   ?restriction3 owl:hasValue ?gene_record.
24.   ?gene_record rdfs:label ?gene_name.
25.   ?process rdfs:label ?process_name.
26. }
```

5 irudia: Neurocommons sistemari SPARQL erabiliz egindako galdera. Erantzunak hainbat jatorriko informazioa elkarlotzen du, eta beharrezko informazio zehatza soilik ematen du, gehiegizko informazioa baztertuz. [<http://bib.oxfordjournals.org/content/10/2/193/F6.expansion.html>]

## 2.3.- HEDABIDEAK HEDATU

NYTk<sup>3</sup> (New York Times) eta BBCk<sup>4</sup>, sona handiko hedabideek, datu estekatuak bi helburuekin erabiltzen dituzte. Alde batetik, irakurleei informazioa igorri aurretiko datuen barne kudeaketa hobetu dute<sup>5</sup>. Bestetik, haien datuak ireki dituzte<sup>6</sup>, irakurleek datuok erraz erabili ditzaten, beste datuekin uztartu ditzaten edo datuen inguruan aplikazioak eraiki ditzaten<sup>7</sup>. NYTk, adibidez, pertsona ospetsuak zein unibertsitateetara joan ziren aurkitzeko bilatzailea<sup>8</sup> garatu du (6 irudia). NYTk argitaratutako datu estekatuen inguruan oso aplikazio interesgarriak sortuko dira datozen urteetan, egunkari horrek gordetzen duen informazioaren kalitatea eta garrantzia kontuan hartuta.

The screenshot shows the 'Linked Open Data' application interface. At the top, it says 'The New York Times' and 'Linked Open Data BETA'. There is a 'View Application Source' link. The main section is titled 'Alumni In The News' and includes a search bar where 'Harvard University' is entered. Below the search bar, there are two profiles: Ed Wilson (born June 10, 1929) and Franklin Delano Roosevelt (born January 30, 1882, died April 12, 1945). Each profile has a list of related articles from The New York Times. On the left side, there is a network diagram with nodes representing various data sources like LIBRIS, Wikicommons, Open Calais, DBpedia, and KEGG.

6 irudia: NYTk eskaintzen duen zerbitzua, pertsona famatuen jatorrizko unibertsitateak aurkitzeko. Bilaketaren emaitzetan pertsona bakoitza aipatzen duten albisteak ere agertzen dira.  
[<http://data.nytimes.com/schools/schools.html>]

3 <http://data.nytimes.com/>

4 <http://www.bbc.co.uk/>

5 [http://www.bbc.co.uk/blogs/bbcinternet/2010/07/the\\_world\\_cup\\_and\\_a\\_call\\_to\\_ac.html](http://www.bbc.co.uk/blogs/bbcinternet/2010/07/the_world_cup_and_a_call_to_ac.html)

6 <http://www.bbc.co.uk/nature/feedsanddata>

7 <http://open.blogs.nytimes.com/2010/03/30/build-your-own-nyt-linked-data-application/>

8 <http://data.nytimes.com/schools/schools.html>

### **3.- DATUEN AMARAUNAREN ETORKIZUNA: WEB SEMANTIKOA**

Datuen amarauna beste proiektu orokorrigo baten lehen pausua da: web semantikoa<sup>9</sup>. Web semantikoaren muina ordenagailuak informazioa "ulertzeko" gai izatea da, erabiltzaileak web sareko informazioaren kudeaketan denbora aurrez dezaten.

Datuen amarauna web semantikoa eraikitzeke oinarri paregabea da, datuen amaraunean datuak RDF lengoia adierazten baitira. Izan ere, RDF lengoia semantikoa da (Datuen esanahia deskribatzen du), eta web semantikorantz lehenengo pausua ahalik eta informazio gehien ordenagailuek kudea dezaketengoa semantikoan adieraztea da. Hurrengo pausua datu horien gainean "ezagutza" sortzea da, ordenagailuek datu horiekin logikoki arrazontzeko gai izan daitezengoa, ordenagailuek informazioa ulertuko balute bezala erabiltzaileontzat lan egin dezaten.

Baliteke web semantikoa inoiz ez ezartzea, baina web semantikora bidean oso teknologia erabilgarriak sortuko ditugu. Datuen amarauna adibide deigarria da, artikulu honetan aurkeztutako aplikazioek erakusten dutenez.

### **ERREFERENTZIAK**

[1] Bizer C, Heath T, Berners-Lee T (2009). Linked Data - The Story So Far. *IJSWIS* 5(3), 1-22

[2] Ruttenberg A, Rees J, Samwald M, Marshall M (2009). Life sciences on the Semantic Web: the Neurocommons and beyond. *Brief. Bioinformatics* 10, 193-204

---

9 <http://www.w3.org/standards/semanticweb/>