



Treball Fi de Carrera

**ENGINYERIA TÈCNICA EN
INFORMÀTICA DE SISTEMES**

**Facultat de Matemàtiques
Universitat de Barcelona**

**ANÀLISI D'ARQUITECTURES BI PER AL
TRACTAMENT DE GRANS VOLUMS DE
DADES BIG DATA I CLOUD STORAGE**

Noël Torres-Caballé

Director: Enric Biosca Trias
Realitzat a: Departament de Matemàtica
Aplicada i Anàlisi. UB

Barcelona, 17 de Gener de 2012

A qui més ha lluitat per aconseguir aquesta fita, la iaia.

Índex

1.	Introducció	4
1.1	Àmbit del projecte.....	4
1.2	Motivació.....	5
1.3	Objectius generals.....	10
1.4	Objectius específics	11
1.5	Organització de la memòria	12
2.	Antecedents. Les dades d'ahir i d'avui.....	14
2.1	El Problema. La ubiqüitat de la informació	16
2.2	Desmitificant el terme “Big data”	18
2.3	Reinventant el consum de la informació.....	21
2.3.1	El processament i l'anàlisi de dades tradicional.....	21
2.3.2	La naturalesa canviant del “Big data”	23
3	Anàlisi. Negoci i Tecnologia.....	25
3.1	Anàlisi de mercat.....	26
3.2	Comparativa i selecció de les eines.....	31
3.3	Orientació de Negoci.....	37
3.4	Escenaris reals d'Implantació per a una solució “Big data”	37
4	Disseny del Sistema.....	40
4.1	Decisions de disseny.....	40
4.2	Arquitectura	43
4.2.1	Capa d'Emmagatzematge.....	44
4.2.2	Capa d'Aplicació	47
4.2.3	Capa de Presentació.....	52
5	Implementació del Sistema.....	54
5.1	Requeriments principals	54
5.2	Construcció del laboratori de proves	57
5.2.1	Instal·lació de CentOS 6.2 com a Sistema Operatiu Virtual	57
5.2.2	Instal·lació del programari de suport.....	62
5.2.3	Instal·lació dels components de la plataforma “Big data”	63

6	Estudi de Costos i Viabilitat.....	70
6.1	Anàlisi del temps de realització del projecte	70
6.2	Valoració del cost econòmic del projecte	70
6.3	Retorn d'Inversió.....	72
7	Conclusió	73
8	Referències Bibliogràfiques.....	75
I.	Annex. Manual d'Usuari de l'entorn de laboratori	76

1. Introducció

1.1 Àmbit del projecte

L'àmbit del projecte queda enquadrat en l'estudi - en base al context econòmic i empresarial actual – disseny i implementació d'un ecosistema tecnològic que faci possible la implantació de sistemes d'informació capaços d'analitzar les diferents fonts d'informació que a dia d'avui el món empresarial es veu en la necessitat de gestionar per tal de: **Conèixer, Predir i Actuar** envers els seus objectius empresarials i capacitat de creixement en el seu mercat.

Amb el propòsit d'exposar i contextualitzar cadascuna de les necessitats i decisions preses per a la consecució de les fites, aquest estudi farà ús de paraules clau associades a l'àmbit d'estudi mirant de mantenir la objectivitat i equidistància entre les dues visions des de les que es pot estudiar aquesta nova àrea de coneixement. La visió més comercial que orienta i exposa els requeriments de negoci que estan experimentant les empreses i la visió tecnològica que representa el que – a la practica – es possible implementar actualment per oferir garanties que donin resposta a les necessitats més immediates de les empreses que vulguin adoptar aquests mecanismes.

1.2 Motivació

El terme “Big data” – que sens cap dubte segueix oferint definicions diferents segons a les fonts a què ens adrecem – no forma part ja només d’un àmbit estrictament tecnològic, encara que inicialment fou encunyat per referir-se a tots aquells mecanismes necessaris per al tractament de l’actual volum d’informació que les corporacions tenen al seu abast.

Avui però, es entès com una prioritat per aquelles empreses que desitgen ampliar el seu coneixement més enllà de l’anàlisi clàssic de la informació coberta en gran part per la disciplina del *Business Intelligence*¹ i que esperen poder obtenir coneixement de l’**evolució** del seu negoci.

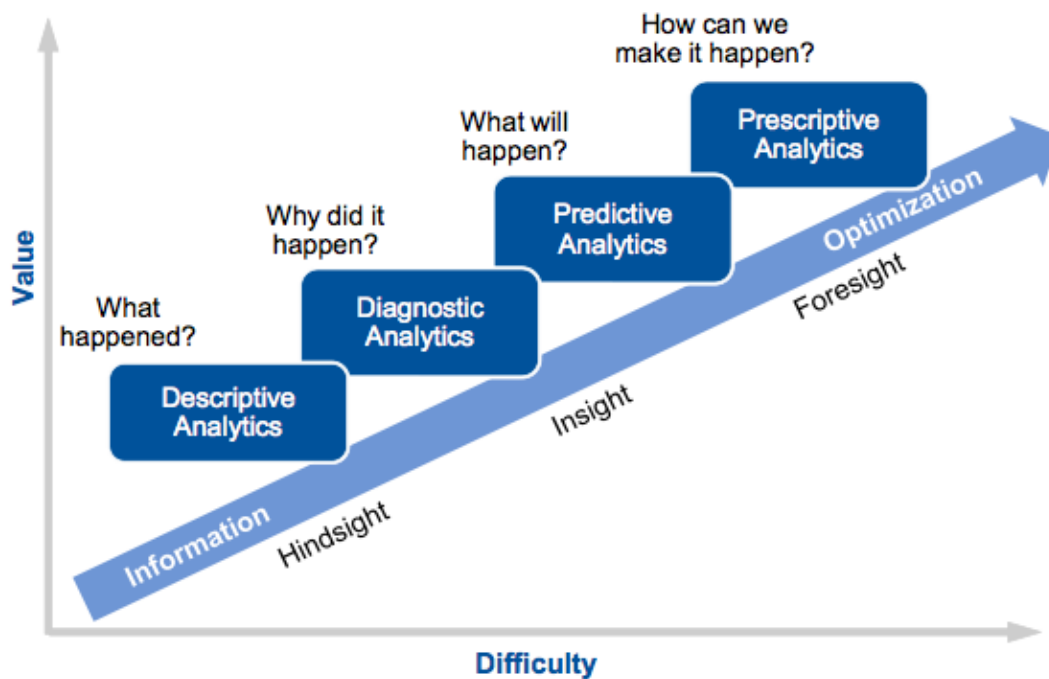


Figura 1. Model d'Ascendència Analítica. Font Gartner 2012

En els models d'anàlisi clàssics¹, les eines podien oferir coneixement al voltant de principalment, dues qüestions: **Què ha passat?** i **Quan ha passat?** (els dos primers estadis del model)

¹ [Business Intelligence](#). Font Wikipedia (2012)

Trobar resposta a aquestes preguntes oferia al negoci la capacitat de **qualificar** els resultats obtinguts en funció de decisions dutes a terme en el passat i **quantificar** els resultats aconseguits.

El nou paradigma “**Big data**”, ens ofereix mecanismes per **qualificar** i **quantificar** la informació en base als dos segons estadis del model: **Què passarà?** i **Cóm pot passar?**

Un recent estudi realitzat pel fabricant *IBM* en col·laboració amb la *Saïd Business School at the University of Oxford*² aporta la visió de l’actual empresariat sobre la manera en que pot influir l’arribada d’aquest nou mètode per a l’anàlisi de la informació:

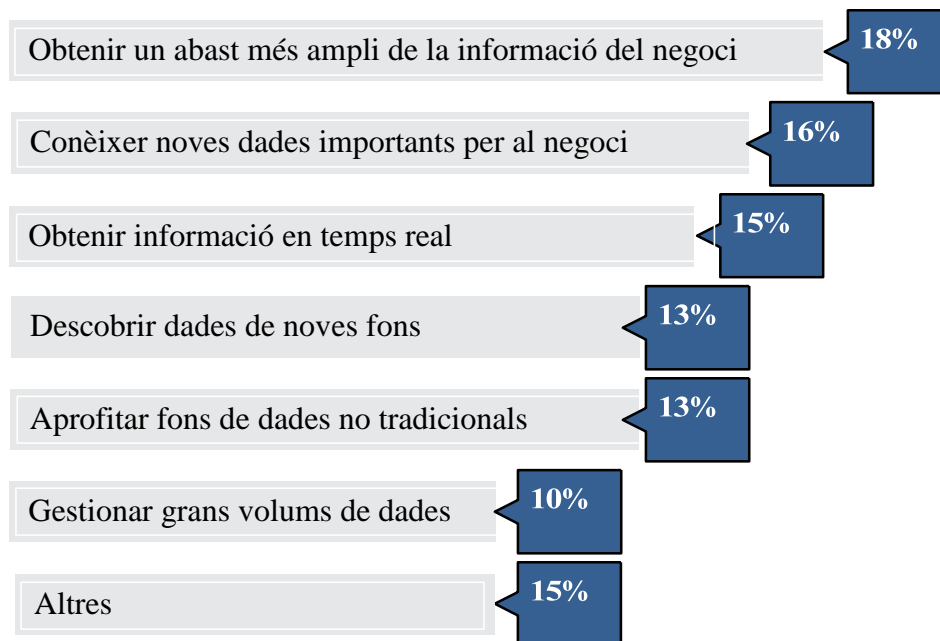


Figura 2. Visió sobre “Big Data” en perfils d’Alta Direcció

Aquests resultats s’alineen amb una manera útil de caracteritzar les tres dimensions del paradigma “**Big data**”- les tres ‘*V*’³-. **Volume (Volum)**, **Velocity (Rapidesa)** i **Variety (Diversitat)**.

² Analytics: The real-world use of big data. IBM 2012

³ “[3D Data Management: Controlling Data Volume, Velocity, and Variety.](#)” Doug Laney 2001, META Group – ara Gartner, Inc -.

Aquestes tres dimensions poden ser visualitzades de forma unificada en la següent figura.

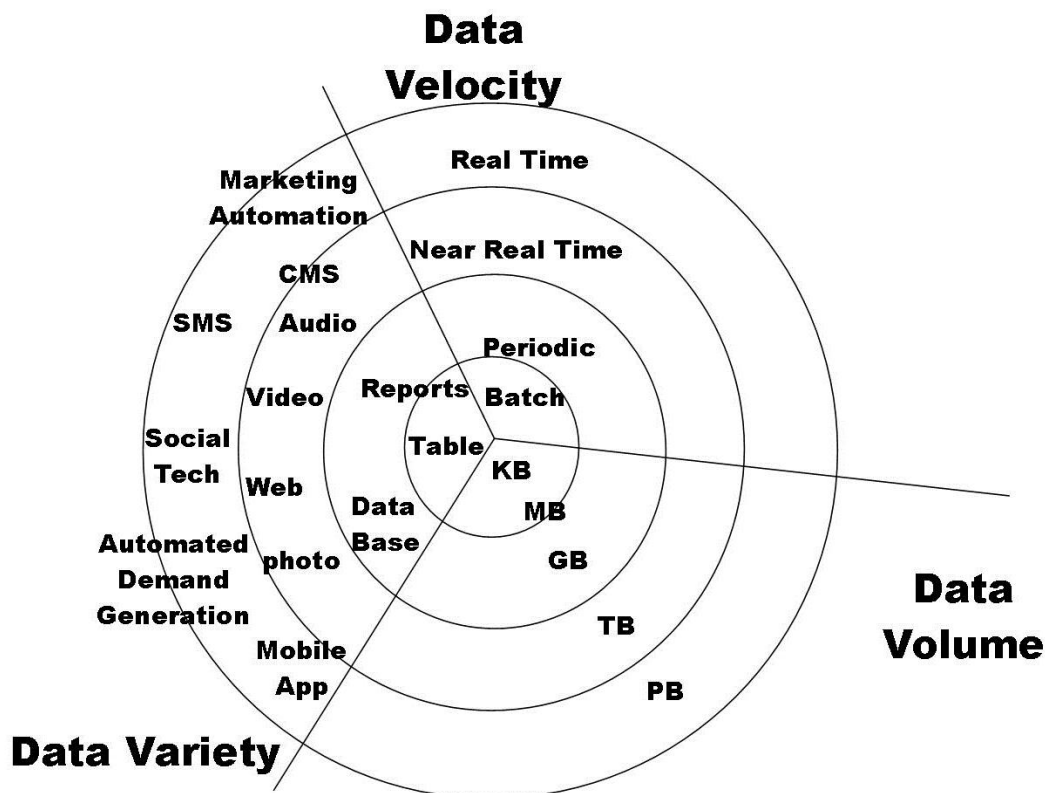


Figura 3. Dimensions d'Anàlisi del paradigma "Big Data"

Data Volume (Volum): La quantitat de dades que gestiona el negoci. Aquesta és, probablement la característica que més comunament s'associa al paradigma del "Big data". El volum de dades generades actualment al voltant d'un sistema d'informació incrementa de forma exponencial i esdevé crític el seu tractament en espais de temps finits.

Data Velocity (Rapidesa): Aquest eix d'anàlisi qualifica la rapidesa en que la informació és generada, processada i analitzada. El període que transcorre entre la captura de la informació i la possibilitat de ser explotada pel negoci – latència – pot ser d'importància crítica per a la presa de decisions.

Data Variety (Diversitat): Les fons d'origen d'informació de la qual les empreses nodreixen el seu coneixement han deixat de ser els orígens tradicionals com els fitxers, i les bases de dades relacionals per passar a procedir de – literalment – infinits sistemes de generació de dades. Poder categoritzar i unificar la informació independentment del seu origen permetrà generar un coneixement més ampli d'aquells events que incideixen sobre el negoci.

Finalment caldrà, amb el propòsit de contextualitzar aquesta categorització, garantir que les dades que afectin al negoci, hagin estat prèviament filtrades i per tant siguin **veraces**⁴.

⁴ IBM completa els eixos d'anàlisi definint la quarta 'V'. Prové de la paraula anglosaxona Veracity.

Més enllà del fort escepticisme que genera aquest nou paradigma sobre si els mecanismes d'abordar la problemàtica actual en la comprensió i l'aprofitament de la informació es l'adequat.

Si bé les xifres parlen per si soles (figura 4), l'estudi que aquí es presenta pretén oposar-se a l'argument per al que el terme “Big data” s'ha encunyat únicament dins una terminologia estrictament de màrqueting per a generar noves oportunitats de negoci⁵.

Year	Enterprise Software Spending for Specified Sub-Markets	Forecast: Social Media Revenue, Worldwide, 2011-2016	Big Data IT Services Spending	Total
2011	2,565	76	24,407	27,047
2012	2,918	1,384	23,476	27,778
2013	3,516	1,812	28,578	33,906
2014	4,240	2,827	37,404	44,472
2015	5,207	3,615	36,189	45,010
2016	6,461	4,411	43,713	54,586

Note: Accuracy is to the nearest \$1 million and is derived from percentage-based algorithms. Research data does not provide total accuracy to the nearest \$1 million. Some rounding errors apply.

Figura 4. Total de despesa en TI generada per l'adopció de Big Data (en milions de dòlars).
Font Gartner Octubre 2012.

Així doncs, i en base als plantejaments disposats en el punt 1.2, l'estudi que aquí es presenta pretén exemplificar una solució al voltant de 4 punts principals:

- 1 De quina manera és possible, amb la tecnologia actual, construir una arquitectura que segueixi el paradigma “Big data”,
- 2 Quins són els condicionants – teòrics i pràctics – que vénen imposats a la adopció d'aquesta tecnologia,
- 3 Com quantificar els costos associats a una iniciativa d'aquest tipus,
- 4 En quin context pot generar un retorn d'inversió satisfactori.

⁵ [McKinsey. Harvard Business Review.](#)

1.3 Objectius generals

L'objectiu d'aquest projecte final de l'ensenyament de Enginyeria Informàtica Tècnica de Sistemes és consolidar els coneixements adquirits durant la carrera i enfocar les capacitats adquirides al llarg de l'ensenyament en l'estudi, anàlisi i desenvolupament d'aquells mecanismes que permeten la implementació d'un sistema d'informació de tipus "Big data".

Objectius Tècnics

Els objectius tècnics estaran enfocats en l'obtenció del coneixement necessari sobre el funcionament de sistemes de fitxers d'emmagatzematge distribuït i algorismes avançats de computació paral·lela que permeten el processament d'informació de manera més ràpida i eficaç.

El coneixement de llenguatge SQL és exigint per treballar amb capes de abstracció que potencien les capacitats de la plataforma. De la mateixa manera es pretén assolir un coneixement avançat de la gestió de recursos de plataformes basades en Linux.

Objectius Personals

Pel fet de ser una tecnologia estretament vinculada a la meua activitat professional, el paradigma del "Big data" i els seus mecanismes, són la evolució natural del Business Intelligence, disciplina en la que porto aplicant el meu coneixement des de fa prop de cinc anys i que està evolucionant de manera molt ràpida per la necessitat d'abastar noves arquitectures capaces de processar volums més grans d'informació en temps real.

Aquest valor afegit és d'una importància estratègica a l'hora d'oferir als clients noves solucions que millorin el seu coneixement del negoci i completar els meus estudis acadèmics amb un treball al voltant d'aquesta àrea marcaran alhora el final d'una etapa – l'acadèmica – com l'inici d'una de nova – la professional – en la que caldrà tenir el coneixement teòric i pràctic per poder oferir solucions en àmbits reals.

1.4 Objectius específics

Particularment els objectius a assolir amb aquest treball d'investigació són els següents.

1. Contextualitzar l'actual estat de l'art en quan a com es gestiona actualment les dades en els sistemes d'informació de les empreses.
2. Analitzar les mancances actuals de les empreses a l'hora de tenir una visió 360° del seu negoci.
3. Proposar les millors eines per cobrir les mancances presentades i plantejar les solucions.
4. Dissenyar una arquitectura de tipus "Big data" capaç de oferir una solució pràctica als problemes prèviament exposats.
5. Implementar el sistema amb les eines proposades.

1.5 Organització de la memòria

La memòria que el lector té a la seva disposició conté els següents capítols, el contingut dels quals passa a descriure's a continuació.

1. Introducció

En aquest punt es presenta l'objecte d'estudi de la memòria, els objectius traçats així com una breu introducció al respecte per situar al lector en el context.

2. Antecedents. Les dades d'ahir i d'avui

El segon punt ubica al lector en el moment actual envers l'objecte d'estudi. I planteja les preguntes clau:

1. **d'Ón** sorgeix la necessitat de incorporar nous mecanismes de tractament de informació? ,
2. **Cóm** es possible abordar les problemàtiques sorgides arrel d'aquests nous mecanismes? ,
3. **Què** aporten a les noves tendències d'anàlisi de les dades? ,
4. **Perquè** l'actual esdevindrà un canvi irreversible en la manera de consumir la informació?

3. Anàlisi. Negoci i Tecnologia

Una vegada establert el terreny i l'àmbit del projecte es necessari cercar els mecanismes que permetran oferir solucions a les qüestions proposades. Aquest capítol realitza aquesta tasca.

4. Disseny del Sistema

En un estrat més baix, aquest capítol descriu formalment les necessitats tecnològiques que caldran per generar un laboratori de proves amb l'objectiu de simular el comportament d'una arquitectura que respongui a les especificacions del paradigma del **"Big data"**.

5. Implementació del Sistema

Seguint amb el procés de construcció de l'entorn, i una vegada definit el seu disseny; tots els procediments per a la construcció d'aquest entorn quedaran descrits en aquest capítol.

6. Estudi de Costos i Viabilitat

Un darrer capítol té el propòsit de unir novament els dos móns – empresarial i de negoci – i tancar el cercle aportant conclusions en termes d'esforç per a la consecució dels objectius.

En aquest mateix capítol es realitza una proposta de viabilitat amb vistes a realitzar una hipotètica incorporació del projecte sobre un entorn productiu en un àmbit empresarial.

7. Conclusió

El darrer punt oferirà una conclusió en relació al projecte, quantificarà les fites assolides i contextualitzarà l'esforç realitzat dins tot el desenvolupat al llarg dels estudis de Enginyeria Tècnica en Informàtica de Sistemes.

2. Antecedents. Les dades d'ahir i d'avui.

A finals del segle XIX, quan el govern dels EEUU va decidir realitzar un cens nacional, els empleats de l'Oficina de Cens, varen enfrontar-se a la difícil tasca de comptabilitzar les més de 60 milions de persones del país mitjançant un dolorós recompte manual de les dades disposades en les partides de naixement que les empreses remetien al govern, a les fitxes personals de cens.

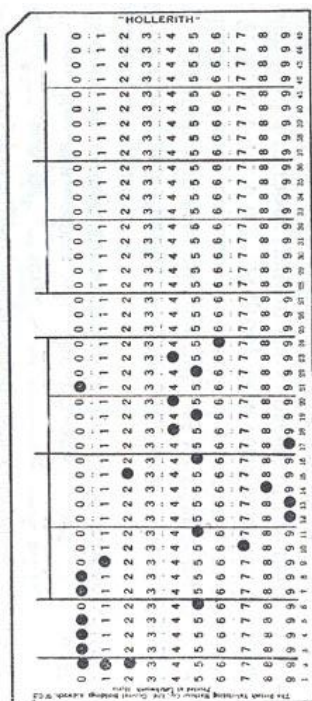


Figura 5. El sistema de Hollerith emmagatzemava dades en forma de forats rodons en una targeta de 45 columnes.

Horroritzat davant la perspectiva, Herman Hollerith⁶ aportà al problema la seva màquina: el **pantògraf tabular**, el funcionament del qual es basava en les màquines perforadores dels conductors del tren als bitllets de viatge per tal d'evitar el frau.

La idea de Hollerith fou l'ús d'una targeta perforada que contenia les dades del Cens dels enquestats i que la màquina elèctrica tabulació podia llegir en pocs segons.

La màquina de Hollerith va processar amb èxit ni més ni menys que 62.622.250 de persones als EEUU, estalviant a l'Oficina del Cens uns \$5 milions de dòlars i reduint el temps de finalització del Cens d'una mica menys de deu anys a prop de 24 mesos.

Per primera vegada s'havia donat una solució a una de les problemàtiques associades al terme "Big data", el **Volum**.

⁶ Martin, T.C., "Counting a Nation by Electricity", *The Electrical Engineer*, New York, November 11, 1891

Un segle més tard, l'any 1997, durant una de les reunions de seguiment de projectes a les oficines centrals de la NASA a la ciutat de Washington D.C, un equip d'investigadors de l'agència aeroespacial tenien com objectiu una fita compartida per diferents equips d'investigació aeronàutica repartits pel món per aquells anys: analitzar e interpretar les dades sorgides de simulacions de motors i dispositius electrònics de mecanismes complexos.

En particular, per a l'equip dels doctors Michael Cox i David Ellsworth, es tractava de recollir tot el conjunt d'indicadors i paràmetres generats per corrents d'aire provocades al voltant d'una nau aeroespacial.

Les super-computadores no eren capaces de processar ni visualitzar tota la quantitat d'informació recollida i oferir-la als investigadors per a un posterior anàlisi i comprensió de la mateixa.

A l'article *Application-controlled demand paging for out-of-core visualization*⁷ Cox i Ellsworth exposaven el següent:

“ [D]ata sets are generally quite large, taxing the capacities of main memory, local disk, and even remote disk. We call this the problem of big data. ”

Aquest es el primer fet documentat on es fa ús del terme “Big data” per tal de referir-se a la problemàtica en el tractament de les dades i exposa els altres dos condicionants que afecten a la resolució d'aquesta problemàtica: **Velocitat** i **Variació**.

⁷ “Proceedings of the 8th conference on Visualization '97”, pàgines 235 i posteriors. IEEE Computer Society Press Los Alamitos, CA, USA ©1997. ISBN:1-58113-011-2

D'ençà d'aquella època – poc abans del canvi de segle – el problema plantejat per uns investigadors ha traspassat els entorns asèptics i controlats dels laboratoris d'alta tecnologia, per convertir-se *de facto* en un repte per aquelles entitats – públiques i privades – que desitgin determinar amb la major exactitud possible quin espai ocupen en el seu àmbit de negoci i conèixer i analitzar tota aquella informació per tal de poder avançar amb més eficiència cap a la consecució dels seus objectius.

2.1 El Problema. La ubiqüitat de la informació

La 'sobrecàrrega d'informació' s'ha convertit en un dels mantres sovint més repetits del nostre temps.

Els llibres s'estan digitalitzant, diaris i revistes ara representen només una fracció dels mitjans de comunicació d'avui en dia, que alhora esdevenen inundats per onades de contribucions deslocalitzades (*tweets*) i publicacions personalitzades (*posts*) a blocs, contingut multimèdia (*video* i *audio*) dona ressò a la informació alhora que els dispositius que fem servir per mantenir-nos al dia amb aquest voràgine, ens permeten fer-nos-en ressò d'una manera ràpida, en qualsevol format i des de qualsevol lloc.

El problema plantejat doncs, no es ja la obtenció, tractament i consum de la informació que una empresa es capaç de generar en si mateixa – i que pot controlar mitjançant els mecanismes àmpliament reconeguts per a realitzar un anàlisi d'aquestes dades – sinó d'aquella informació que sense formar part intrínseca del negoci, pot aportar un valor afegit determinant del que extrapolar indicadors de negoci externs coneixent la manera en què és visualitzada i considerada la empresa a través de tots aquells canals d'informació aliens a la organització.

I doncs, d'on prové aquesta informació?

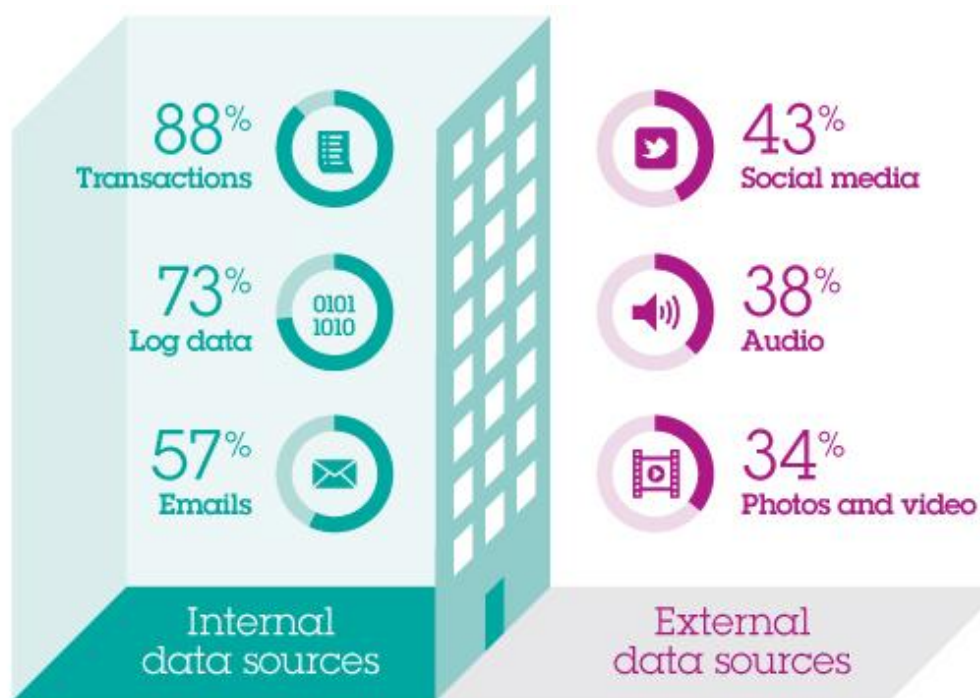


Figura 5. d'On extreuen informació les empreses. IBM, Octubre 2012.

Segons l'estudi de IBM en l'adopció de noves tendències de recollida d'informació, a dia d'avui i de mitjana el 70 % de la informació que manipulen i amb la que es nodreixen les empreses procedeix dels seus propis sistemes d'informació i només en una proporció del 30% aprofiten dades que, si bé són generades per tercers, contenen informació relativa al seu negoci.

Aquest subconjunt d'informació, que actualment no és considerat pels sistemes d'informació de les empreses, conforma el 90% de la informació que, en l'argot comú, s'inclou dins de l'àmbit de "Big data":

“ Big Data, es un terme genèric que pretén donar cabuda a tot aquell conjunt d'informació des-estructurada o semi-estructurada (en endavant: multi estructurada) que pel seu volum o complexitat esdevindria molt costós poder manipular en un sistema clàssic de gestió de la informació -- bé sigui en termes econòmics o de temps -- per analitzar-la. ”

2.2 Desmitificant el terme “Big data”

Per totes aquestes raons, i amb el propòsit de enfocar objectivament les fites que volem assolir quan ens proposem abordar la implantació dels mecanismes que contemplin l'ús d'aquest nou paradigma, serà necessari alliberar-nos de idees preconcebudes al voltant d'aquest terme que poden limitar l'abast i comprensió dels mecanismes que el conformen.

1. “Big data” no fa únicament referència a quantitats ingents de dades.

Primer, la quantitat de la informació és només una de les consideracions d'aquest paradigma i, ampliant l'argumentació en aquest punt de la memòria, és probablement el menys important de tots ells.

Les màquines o instal·lacions modernes, com ara automòbils, trens, centrals nuclears o avions han augmentant constantment el nombre de sensors que recol·lecten informació per al control del seu funcionament. És molt comú tenir milers o fins i tot centenars de milers de sensors per tal de recollir i constatar el correcte acompliment de les activitats d'un sistema.

Una aeronau es compon de aproximadament cent mil sensors que durant una hora de vol recullen tot tipus de dades: la velocitat de l'aire sobre cada part de l'armadura d'avió, la quantitat de diòxid de carboni a cada secció de la cabina, la pressió atmosfèrica, etc... on cada sensor és en realitat un dispositiu independent amb les seves pròpies característiques físiques.

L'interès real de tot aquest muntatge es extreure informació a partir de les **combinacions de lectures dels diferents sensors** (com ara de quina manera el diòxid de carboni es combina amb temperatura de la cabina o bé com evitar les turbulències contrastant la velocitat d'aire amb la pressió que exerceix sobre la estructura).

Amb tants sensors (**Diversitat**) emeten informació a tanta velocitat (**Rapidesa**) les combinacions esdevenen increïblement complexes i pot variar en funció de la tolerància d'error segons les característiques dels dispositius individuals.

Així doncs, les dades proporcionades per un centenar de milers de sensors de que es compona una aeronau pot encabir-se dins els paràmetres de “Big data”.

No obstant això, la grandària del conjunt (**Volum**) de dades no és tan gran com es podria esperar. Amb uns càlculs ràpids podem exemplificar-ho senzillament. Per a un conjunt de 100.000 sensors, on cadascun dels quals produeixin vuit bytes d'informació per segon generaria menys de 3 GB de dades en una hora de vol

$$100.000 \text{ sensors} \times 60 \text{ minuts} \times 60 \text{ segons} \times 8 \text{ bytes} = 2.8 \text{ GB.}$$

2. El paradigma “Big data” està orientat a l'anàlisi, no a l'emmagatzematge

Segon, les companyies no necessiten aquelles dades a partir de les quals ja n'han extrapolat resultats, generat coneixement i pres decisions.

Emmagatzemar, consultar y recolzar les accions dues a terme sobre el negoci en base a la **informació històrica** d'una companyia es l'objecte d'estudi de les metodologies clàssiques proposades per Inmon, Kimball o més recentment Linstedt ⁸ en la implantació de sistemes Data warehouse⁹ i correspon a l'àmbit de l'àrea de Business Intelligence, i escapa del propòsit final d'una arquitectura basada en el paradigma del “Big Data”.

Tota aquella informació que pogués ser emmagatzemada en un sistema d'aquestes característiques ha recorregut un procés previ de transformació i homogeneïtzació per ajustar-se a les necessitat d'anàlisi del negoci.

⁸ Inmon, Bill (1992). *Building the Data Warehouse*. Wiley. [ISBN 0-471-56960-7](#). Kimball, Ralph (1996). *The Data Warehouse Toolkit*. Wiley. [ISBN 0-471-15337-0](#). Linstedt, Graziano, Hultgren. *The Business of Data Vault Modeling Second Edition (2010)* Dan linstedt, [ISBN 978-1-4357-1914-9](#)

⁹ [Data warehouse](#). Font Viquipèdia

3. El paradigma “Big data” no aplica únicament a empreses grans

Tercer, la decisió d’incorporar una arquitectura d’aquestes característiques en cap cas depèn de la dimensió o magnitud de la companyia que desitja fer-ne ús.

Per a una petita o mitjana empresa amb un reduït equip de personal, que utilitza un sistema de gestió de la informació transaccional¹⁰, manté un negoci online i desitja aprofitar la informació que generen els sistemes de logs de la seva web, aconsegueix el perfil per iniciar una proposta sota una arquitectura “Big data”.

Òbviament aquesta decisió caldrà ser suportada per una plataforma tecnològica adient, si bé els condicionants del negoci per a l’adopció d’aquesta filosofia d’anàlisi estarien plenament coberts.

4. El paradigma “Big data” significa dades des estructurades

El terme “des estructurat” no és precís i no engloba les diferents i variades estructures típicament associades amb els tipus de dades que comunament s’associen al “Big Data”.

Per aquest motiu, el terme que millor s’adequa per aquests tipus de dades és “multi estructurat” al poder incloure cadenes de text, documents de tot tipus, fitxers d’àudio i vídeo, metadada, pàgines web, missatges de correu electrònic, xarxes socials, formularis, etc...

Cadascuna de les tipologies presentades comparteixen totes elles un tret comú:

“*L’esquema en què es basen les dades no és conegut o bé no està definit en el moment que les dades es capturen i s’emmagatzemen. El model de dades que les gestionarà s’aplicarà en el moment en què s’utilitzin.*”

¹⁰ [OLTP](#). Font Viquipèdia

2.3 Reinventant el consum de la informació

Fins a dia d'avui les empreses han consumit la informació que utilitzen en el seu negoci d'una mateixa manera.

La nova situació a la que les corporacions s'enfronten per tal de absorbir tota la informació generada al seu voltant, han obligat a canviar la manera en que processen les dades. Aquest capítol posa èmfasi en descriure aquest punt per tal que el lector prengui consciència de les mesures necessàries que, a la pràctica, caldran dur-se a terme per a possibilitar un tractament eficient de tota la informació.

2.3.1 El processament i l'anàlisi de dades tradicional

Tradicionalment, el tractament de dades amb finalitats analítiques ha seguit un model eminentment estàtic.

És a dir, al llarg del procés habitual del negoci, les empreses generen diferents quantitats d'informació sobre models de dades estables mitjançant aplicacions empresarials com CRM's, ERP's i/o sistemes de control financer. Les eines d'integració de dades s'empren per tal d'extreure, transformar i carregar la informació generada per les aplicacions empresarials i bases de dades transaccionals sobre una "zona d'espera" on té lloc un control de la qualitat de les dades i la seva normalització i les dades passen a integrar-se en un model perfectament acotat de fileres i columnes.

Les dades modelades, i netejades de inconsistències passen a carregar-se en un magatzem de dades empresarial o data warehouse.

Aquesta rutina acostuma a tenir lloc de forma regular i amb periodicitats establertes, en base a les necessitats de la empresa.

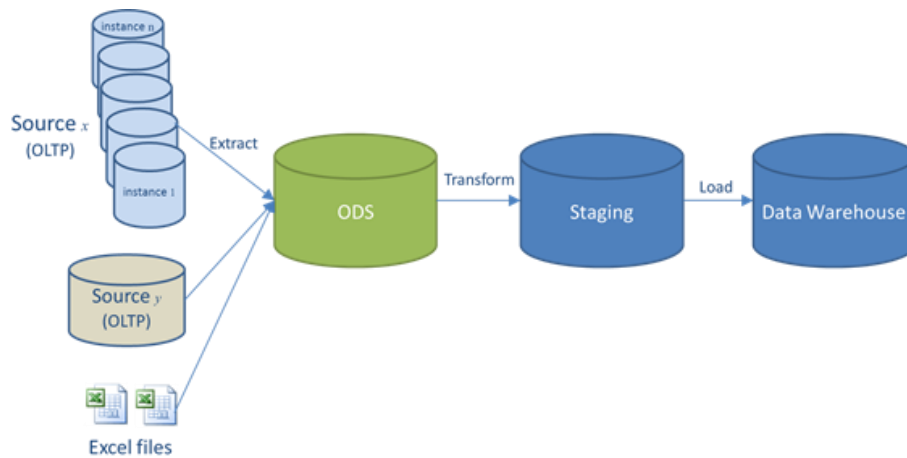


Figura 6 – Model Tradicional de Processament i Anàlisi de la informació.

A partir d'aquest moment, comença la tasca dels administradors i analistes de base de dades per tal de crear i programar informes periòdics que funcionin sobre les dades normalitzades i emmagatzemades al data warehouse, que un cop construïts es distribueixen als usuaris de negoci.

Son també els responsables de crear quadres de comandament i demés eines de visualització per aquells perfils executius i de gestió que necessitin veure la evolució del negoci. Els analistes de negoci, per la seva banda, utilitzen les eines d'anàlisi per realitzar consultes avançades contra el magatzem de dades i els usuaris no experts realitzen la visualització i anàlisi de dades bàsiques limitats pels objectes d'exploració i el front-end de les eines de Business Intelligence.

Els volums de dades en magatzems de dades tradicionals amb prou feines poden superar els centenars de terabytes on l'increment mal dimensionat del volum de la informació pot degradar el rendiment del sistema.

2.3.2 La naturalesa canviant del “Big data”

L'arribada de la web, els dispositius mòbils i la incursió del món digital en el nostre sistema de vida han provocat un canvi fonamental en la naturalesa de les dades.

El “Big data” té qualitats importants i diferenciadores que el desmarquen de la concepció tradicional de les dades corporatives. Ja no s troben centralitzades, responen a estructures constants o són senzillament transformables. Ara més que mai, la informació es troba distribuïda, apareix vagament estructurada (si pot ser adaptada a algun patró específic), i augmenta cada vegada més en tamany.

En termes generals, aquesta informació bé generada per fonts tals com:

- **Xarxes socials i mitjans de comunicació:** Actualment hi ha més de 700 milions d'usuaris a Facebook, 250 milions d'usuaris a Twitter i 156 milions de blocs públics.

Qualsevol actualització d'aquestes fonts crea nous *punts de dades*, multi estructurats també anomenat “*Data exhaust*”¹¹.

- **Dispositius mòbils:** Hi ha més de 5 bilions de telèfons mòbils en ús a tot el món. L'ús d'aquest dispositiu genera dades en tots els sentits: la durada de la trucada o missatge, la ubicació, el destí, etc... Els dispositius mòbils, telèfons intel·ligents i les tauletes en particular, promouen l'ús i faciliten l'accés a les xarxes socials així com altres mecanismes de generació de dades.

- **Les transaccions d'Internet:** Amb milers de milions de compres en línia, les accions comercials tenen lloc contínuament, cada segon. Cada una d'elles genera incomptables punts de dades recopilades pels establiments, entitats financeres, ISPs, agències de crèdit, etc...

¹¹ *Data Exhaust: Aquelles dades i/o informació multi estructurada, producte de les activitats generades per usuaris en entorns on-line.*

- **Els dispositius connectats en xarxa i sensors:** els dispositius electrònics de tot tipus - inclosos els servidors i altre maquinari de TI, mesuradors d'energia intel·ligents i sensors de temperatura - tot crear semi-estructurats les dades de registre que registren cada acció.

DADES TRADICIONALS	BIG DATA
Gigabtes - Terabytes	Petabytes - Exabytes
Centralitzada	Distribuïda
Estructurada	Multi estructurada
Model de dades fix	Esquemes Plans
Interrelacions complexes de la informació	Poc nivell d'interrelació en el procés de les dades

Figura 7 – Dades Tradicionals vs. Big Data

Així doncs aplicar els conceptes tradicionals de gestió de la informació sobre aquests nou tipus de dades, seria inviable degut al temps i l'esforç necessari que suposaria modelar – en un procés iteratiu - un sistema adequat a les característiques de cada nova entrada susceptible de ser analitzada.

Una altre argument que reforça aquesta posició, en aquest en l'àmbit del negoci, són els costos de TI que haurien de dedicar-se en infraestructura per processar quantitats potencialment tant elevades de dades.

“Sota aquestes condicions es fa evident la recerca de noves formes de processar y analitzar les dades. Benvinguts a l'era del “Big data”.”

3 Anàlisi. Negoci i Tecnologia

El Març de 2012, l'estudi “*Worldwide Big Data Technology and Services 2012–2015 Forecast*” de l'empresa IDC¹² va revelar que el creixement estimat del mercat en l'adaptació del paradigma “*Big data*” en l'àmbit empresarial privat, augmentarà de 3.2 bilions de dòlars que va suposar per a l'any 2010 a 16.9 bilions per a l'any 2015.

Per a un negoci en el que les seves expectatives de creixement són d'aquesta magnitud (a raó de més del 40% anual), seria lògic pensar que les propostes per implantar aquestes mecanismes haurien d'estar – literalment – inundant el mercat.

La realitat, però, es que malgrat la constant aparició d'articles, notícies, seminaris o events al voltant d'aquest tema i a dia d'avui, els fabricants continuen tractant amb molta prudència la seva incursió en aquest terreny, a la expectativa de conèixer la necessitat real de les empreses per a la implantació d'aquesta tecnologia en els seus processos de negoci.

Les nascudes en l'era digital, pioneres en els seu àmbit com *Google*, *Facebook* o *Twitter* i que degut al objecte del seu negoci requereixen d'aquesta tecnologia, varen iniciar-se en l'ús de diferents conceptes, idees i projectes provinents del món Open Source per tal de donar solució a la seva problemàtica, que si bé inicialment eren poc madurs, han anat incrementant el seu grau de sofisticació a mida que el dia a dia obligava a cobrir noves demandes dels seus clients.

En el darrer lustre, han estat un bon grup d'empreses noves les que davant la oportunitat de posicionar-se en el mercat en un àmbit tecnològicament innovador, han adoptat també aquestes tecnologies on veuen l'espai per poder competir de tu a tu amb els grans fabricant amb una idea clara: oferir propostes que simplifiquin, optimitzin i modularitzin el conjunt de solucions tecnològiques que conformen una plataforma de tipus “*Big data*”.

¹² International Data Corporation (IDC) es una prescriptora líder en estudis d'intel·ligència de mercat i serveis d'assessorament tecnològic, telecomunicacions i mercats de consum tecnològics.

I tot això, per tal de convèncer a les empreses amb una proposta atractiva per aquelles clients potencials que es troben en la necessitat de treballar en aquesta direcció podent alhora, aprofitar la seva infraestructura tecnològica per alimentar el seus nous sistemes de informació.

Aquest capítol està dedicat a revisar l'estat actual del mercat envers a la tecnologia que actualment cobreix aquest àmbit i presentar la proposta tecnològica escollida per al desenvolupament del projecte i que s'abordarà en els capítols 4 i 5 d'aquest PFC.

3.1 Anàlisi de mercat

Començant per les empreses més conegudes i consolidades fins a les *start-up's* més agosarades, tothom ha decidit – d'una manera o altre – iniciar la seva aposta per accedir al mercat del “Big data”.

Els reptes d'aquesta nova proposta tecnològica i el valor potencial que aporta en la presa de decisions per a les empreses es un element fortament motivador, on els gerents estan constantment a la recerca de l'avantatge competitiu per tal de guanyar una presència més forta en els seus respectius mercats. Amb tant potencial per proporcionar a les empreses amb tècniques d'anàlisi millorades i un processament de més informació, és comprensible l'interès que aquesta nova àrea ha generat tantes expectatives – i publicitat – associades amb la seva implantació.

El que a continuació es presenta, es un anàlisi de l'actual mercat pel que fa a l'adopció del “Big data”.

Per a la elaboració d'aquest anàlisi s'han revisat, compilat i reflectit informacions publicades per entitats prescriptores durant el període (Març 2012 – Desembre 2012) fet que ha facilitat la confrontació de la informació, correcció i ajustos necessaris per obtenir la visió més acurada possible fins el moment.

S'ha avaluat la demanda d'aquesta tecnologia en diferents indústries i s'ha generat un mapa de calor per mirar d'oferir una visió de les necessitats potencials per a cada sector en l'àmbit del "Big data".

QUINES SERAN LES DIMENSIONS DE NEGOCI MÉS DEMANDADES DEL PARADIGMA "BIG DATA" PER A L'EXPLOTACIÓ DE LA INFORMACIÓ CORPORATIVA?

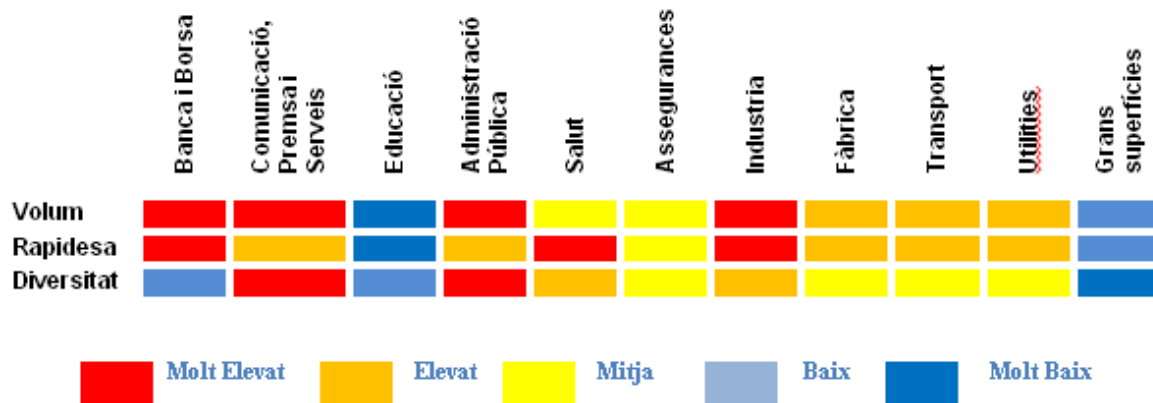


Figura 8. Oportunitats de Negoci estimades en implantacions "Big Data"

A partir de la informació presentada en aquest quadre, es pot concloure que la "V" més perseguida actualment alhora d'adoptar "Big data", es la capacitat de processar més ràpidament la informació. Això es demanat pels sectors de les comunicacions i de la informació (oferir continguts actualitzats en espais de temps més reduïts) o en el terreny de l'administració pública (agilitzar els processos burocràtics actuals); altres sectors centrats en aquesta dimensió d'anàlisi són els de Fàbrica, Transport i Utilities. Aquesta dimensió però, cau en el segon lloc com a **prioritat en la seva implantació**, per darrera de "Volum".

La segona "V" en la que centrem la nostra atenció i que te prevista una implantació més ràpida es la que afecta al tractament de grans quantitats de informació. Es important fer notar que quasi 7 de cada 10 sectors tenen aquesta prioritat com a necessitat a cobrir en els propers anys.

Aquestes dues "V" conformen prop del 60% de les previsions que de forma transversal afecten als diferents sectors d'estudi.

L'estudi també revela que els sectors on el "Big data" te actualment menys capacitat de penetració es el de la **Educació** i el de les **Grans Superfícies**.

Centrant-nos doncs en els diferents sectors econòmics descrits anteriorment, s'ha cercat informació al voltant de quina ha estat fins ara la capacitat d'aquesta tecnologia alhora de fer-se un lloc a les empreses que contribueixen en cada sector.

Per a oferir la informació d'una manera segmentada i el més clara possible, s'han distribuït els percentatges d'adopció en base a una pregunta i els resultats encabits en 5 àmbits diferenciats.

En el següent gràfic es pot revisar el resultat de creuar la informació obtinguda. Val a dir, per ésser rigorós, que la mostra no respon directament a un univers d'empreses, sinó al extrapolat de la recerca, revisió i quantificació de les dades obtingudes. Tot seguit es presenten els resultats obtinguts.

HA INVERTIT LA SEVA COMPANYIA EN ALGUNA TECNOLOGIA ADREÇADA ESPECÍFICAMENT A COBRIR NECESITATS ASSOCIADA AL PARADIGMA "BIG DATA"?

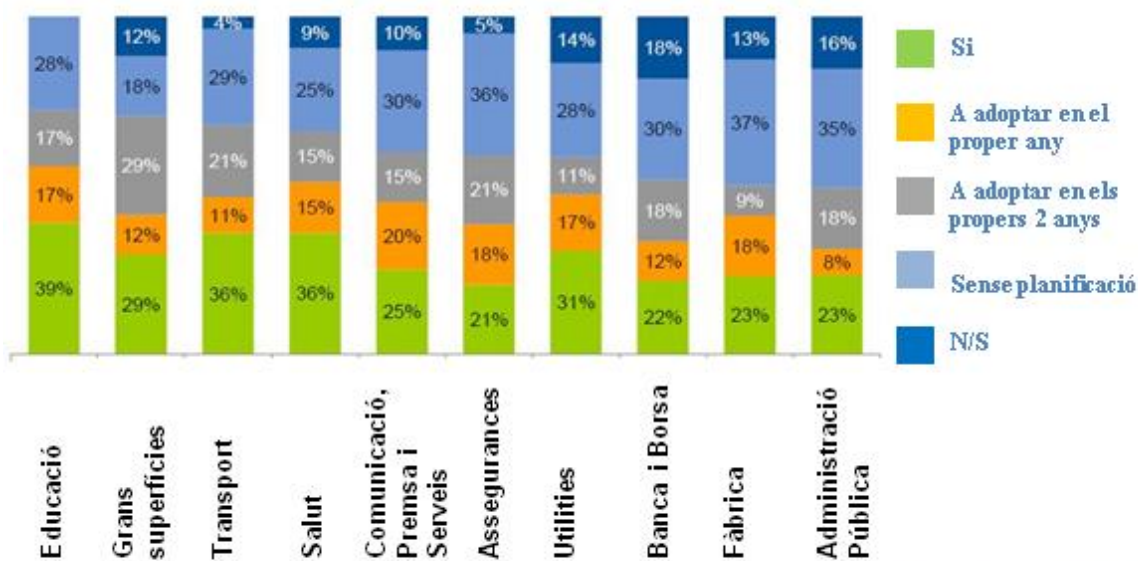


Figura 9. Inversions en "Big Data"

Les dades recollides en aquest àmbit i representades en la figura 9, ens donen una mostra de la implantació actual del “Big data” a nivell sectorial.

Observem que, de mitjana, prop de 4 de cada 10 empreses – amb independència del sector al que pertanyen – tindrien implantat actualment o a un any vista, algun procés associat a oferir solucions a problemàtiques vinculades a temes de “Big data”. Aquest es un percentatge molt elevat tenint present el grau de maduresa de les eines tecnològiques així com l’oferta de productes associats, i es un indicador de l’interès que suscita aquesta proposta per a les empreses en l’intent de demostrar, a la pràctica, la seva efectivitat.

A la banda oposada de la forquilla, aproximadament la mateixa proporció no han pressupostat cap projecte en aquest sentit o bé no s’han assessorat o desconeixen els beneficis que podria reportar al seu negoci l’adopció d’aquestes tècniques.

Com a curiositat, el sector Bancari i de Borsa malgrat ser un consumidor potencial d’aquesta tecnologia, es el que a curt termini en menor mesura (4 de cada 10 empreses) adoptarà aquestes tecnologies. L’alt cost en l’adequació dels sistemes actuals i la criticitat de les dades que aquestes companyies gestionen, es pot revelar com el major handicap per a la inclusió d’aquestes tècniques.

Com a darrer punt i per completar la revisió del marc comercial, hem inclòs una entrada que mirar de retratar – al marge dels àmbits potencials en els que el seu ús seria desitjable – recull les orientacions realistes i pràctiques en les que s’emprarà la tecnologia “Big data”.

Aquesta visió la recull el “*Hyper Cycle Big Data 2012*” que *Gartner Inc*, va publicar durant l’estiu de 2012. L’objectiu es senzill: unificar en una corba de gauss aquells àmbits de potencial aplicació d’aquesta tecnologia i enquadrar-los en cinc estats:

1. Durant l’aparició de la tecnologia,
2. En el moment en que aquesta genera més expectatives,
3. Quan es demostra que la tecnologia no cobreix encara tot el ventall demandat,
4. Durant la corba en que s’estudia detingudament on podria adaptar-se de forma més òptima donades les circumstàncies, i finalment,
5. Quins són els àmbits en que pot esdevenir una tecnologia aplicada i productiva.

En base a aquests espais i incloent la dimensió del temps, es possible establir prediccions sobre de quina manera pot influir – actualment i en el futur pròxim – la tecnologia objecte de d’estudi.

La *figura 10* disposada a continuació, exemplifica aquesta visió.

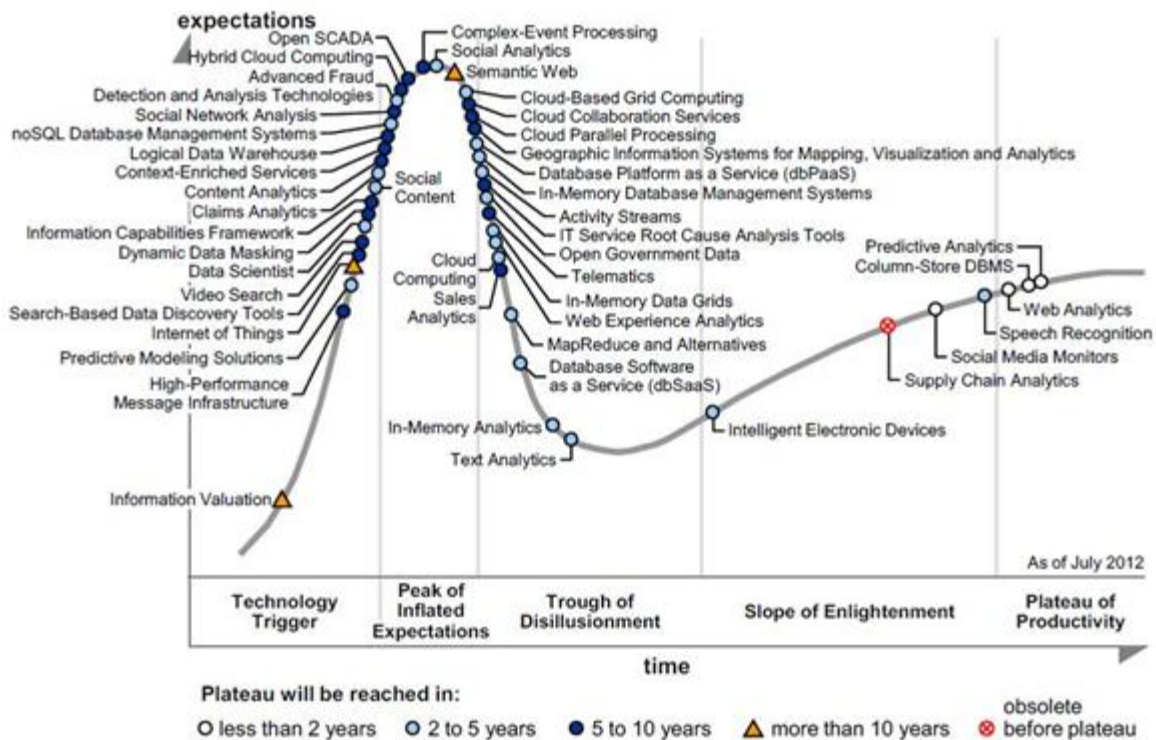


Figura 10. *Hyper Cycle Big Data 2012*. Font *Gartner Inc*. (Juliol 2012)

3.2 Comparativa i selecció de les eines

Per a la consecució dels objectius plantejats en aquest projecte, s'han hagut de revisar les diferents tecnologies que actualment estan presents en el mercat i decidir sobre quina combinació seria més adequat presentar el laboratori en el que mostrar la solució implementada.

Com en qualsevol iniciativa de TI, és molt important per a una organització alinear les necessitats del negoci amb la elecció de la seva tecnologia pel que fa al paradigma “Big data”.

Les necessitats de negocis poden ser infinites: explorar la informació dels clients per tal de predir el seu tipus de despesa, afinar les ofertes amb productes personalitzats o millorar la venda creuada són objectius comuns que exigeixen l'anàlisi constant de la informació per traçar patrons i predir noves situacions.

Aquestes necessitat no s'ajusten a les arquitectures actuals tal i com les entenem, si bé fan ús de les mateixes per adaptar el canvis sorgits arrel d'aquestes noves necessitats. Les organitzacions voldran evitar moure dades no processades prèviament directament a al seu magatzem de dades ja que només disposarien d'això: dades, quan el que cerquen es el nivell superior: la informació.

Les **tres grans arquitectures** que actualment es poden trobar al mercat entorn al “Big data” venen condicionades per principalment 2 factors:

- **El pressupost**, o la capacitat d'una companyia alhora d'invertir en una solució tecnològica,
- **La necessitat**, o l'àmbit d'actuació que la companyia desitja cobrir per millorar les seves capacitats de negoci, i

Un darrer factor que també seria d'interès es la **capacitat tecnològica** de la companyia per autogestionar aquests nous serveis (en el sentit de si escull implantar i formar els seus tècnics o bé decideix subcontractar aquest coneixement) si be, per simplificar, considerarem en aquest punt un model híbrid en el que la companyia externalitza la implantació de la solució en una primera fase i posteriorment assumeix el control de la mateixa formant el seu departament tècnic.

Arquitectura Tipus 1. - Models Propietaris

Els propietaris per aquest tipus de tecnologia requereixen d'una despesa de capital significativa per iniciar el projecte. Objectivament una arquitectura basada en aquest model resulta en un major cost de TCO¹³ donat que els fabricants condicionen les seves propostes a fer ús de llicències de software i components hardware certificats en base a garantir el correcte funcionament de la plataforma.

Sovint els riscos en la adopció de components de codi obert que permetin la construcció de elements tant determinants com la infraestructura de xarxes, la configuració de múltiples servidors, els protocols de seguretat o la personalització de codi es inviable i poc rentable. És molt important valorar aquests punts quan es pensa en construir una plataforma d'aquest tipus. Construir eines personalitzades pot semblar una mesura d'estalvi, però no sempre és la millor opció quan la compra de tecnologia pot salvar incomputables hores de feina, evitar retards o en el pitjor dels casos l'abandonament d'un projecte d'aquesta índole.

En aquest cas trobem fabricants (i les seves respectives solucions empresarials) de la talla de: Teradata, Oracle Exadata, Microsoft, IBM o Informatica.

¹³ Acrònim per "*Total Cost of Ownership*" que fa referència a la despesa global que suposa la implantació d'una solució de TI.

Arquitectura Tipus 2. - “Big Data” Appliances

Incorporar un projecte en aquest format es una manera efectiva d’obtenir solucions que cobreixin les necessitats amb un menor TCO. En general les implementacions d’aquests tipus de propostes permeten abstrure’s del negoci i centrar-se en el modelatge de les dades que volen ser tractades així com en les fases d’anàlisi que es desitgen implementar.

En oposició a la construcció d’un sistema des de zero com es plantejava en el punt anterior, aquesta opció elimina molt temps de configuració de maquinari i components de codi obert.

Hadoop, un sistema de codi obert d’arxius distribuït, ha esdevingut la pedra angular en el tractament d’aplicacions de “Big data”. *Hadoop* no aplica costos de llicències per al maneig i processament de dades; al marge de estar a l’àmbit del codi obert es alhora agnòstic en quan a dades, i permet adaptar-se a aquells negocis que cerquen obtenir coneixements a partir de fonts de dades multi estructurades, com registres web i mitjans socials.

Com a resultat, un gran nombre fabricants ja consolidats en altres àmbits com IBM, EMC, HP o Teradata han adquirit companyies que varen néixer per treballar en aquest àmbit i que disposaven de l’*expertise* tecnològic específic per cobrir aquestes necessitats:

- **IBM** va adquirir Netezza i actualment disposa d’un producte que aporta solució a problemàtiques “Big Data”: Big Insights.
- **EMC** va adquirir Greenplum, el segon és ara el nom de la solució de EMC en aquest terreny.
- **HP** va adquirir Vertica, i el sistema columnar de base de dades és part de la solució “Big Data” d’aquesta companyia.
- **Teradata** va adquirir “Aster Data” per proveir de computerització paral·lelitzada al seus sistemes de “Big Data”.

Finalment, i com hem fet notar anteriorment, bona part de les grans companyies han acabat per adoptar *Hadoop* per donar forma a les seves plataformes “Big Data”.

Alguns d'aquests fabricants que el lector coneixerà són Oracle, amb el seu producte **Exadata** basat en aquesta tecnologia o Microsoft, que finalment va abandonar el seu projecte **Dryad** en favor de *Hadoop* per oferir suport a la tecnologia de fitxers distribuïts.

Finalment, i amb relació directe al que s'exposarà en els capítols 4 i 5 d'aquest estudi, cal mencionar dues empreses que són les responsables de la majoria de contribucions quantitatives i qualitatives al projecte *Apache Hadoop* en què els grans fabricants de software basen les seves solucions: **Hortonworks** i **Cloudera**.

- Els enginyers de Hortonworks treballen actualment en una arquitectura *MapReduce*¹⁴ que podrà ampliar el tamany màxim d'un clúster de *Hadoop* més enllà de la seva limitació actual de 4000 nodes, així com noves capacitats de processament de dades en temps real i altres capacitats avançades d'anàlisi de dades.
- Cloudera, en el seu cas, ha contribuït (i contribueix) a seguir desenvolupant *Apache HBase*, la base de dades no-relacional basada en *Hadoop* que permet tendir a zero els temps d'accés a la informació.

Mitjançant la integració dels components que aquestes empreses proporcionen és possible aconseguir una arquitectura flexible, escalable i compatible quant a dades, poden centrar tota l'atenció en trobar valor a la informació sobre el client de manera més ràpida.

¹⁴ El capítol 4.2 detalla l'objectiu de les tecnologies emprades i el seu ús.

Arquitectura Tipus 3. - Plataformes “Big Data” al ‘núvol’

Una darrera opció: pensar en una arquitectura basada en el ‘núvol’.

Disposar una solució “Big Data” en aquest format permet al client beneficiar-se de l’anomenat “Investigative Data Lab” (IDL). L’IDL es basa en la ubiqüitat a què respon tant la **solució** com a **les dades que es desitgen analitzar** i permet als usuaris de negoci accedir a les dades a la recerca de *senyals* associades a elles en les xarxes socials, les dades pròpies de la companyia i les dades multi estructurades.

Quan els senyals es troben – es descobreixen patrons coincidents – poden ser **socialitzats** a través de les línies de negoci, **consultats** en el núvol (no cal pensar ja en recuperar les dades de còpies de seguretat de l’empresa o ubicacions de xarxa, etc...), **incorporats** i **organitzats** en el magatzem de dades i **compartits** de forma transversal a tota la companyia. Recórrer una gran quantitat d’informació és molt costós en temps: centrar la recerca en l’acompliment de patrons esdevé més àgil i eficaç.

Aquest model encara està poc implantat per la seva complexitat tecnològica i la recança que les empreses encara mostren a l’hora de deslocalitzar la seva informació adduint raons de seguretat i privacitat.

Arribats aquest punt i havent realitzat un acostament quant a l'arquitectura, el projecte es basarà en la construcció d'un "**Big Data**" **Appliance** mitjançant l'ús de les següents eines on, per als capítols 4 i 5, en descriurem amb detall tant el del seu funcionament com la manera en què s'ha integrat en el conjunt de la solució "**Big data**".

- **Hadoop Distributed File System (HDFS)** per l'emmagatzematge de dades.
- **MapReduce** per al processament dels conjunts de dades – data sets – sobre clústers.
- **HBase** per a proveir un ràpid accés I/O sobre dades.
- **Hive / Pig** per a l'execució de consultes tipus SQL sobre els conjunts de dades - data sets – sobre models columnars mitjançant *RCFile*.
- **Sqoop**
- **Flume** per a l'estudi de col·leccions de dades en temps.
- **JDBC** i **ODBC** per possibilitar la interconnectivitat entre SGBDR i
- **Hue** com a interfície d'usuari.

L'arquitectura es completarà amb una sèrie de serveis associats – *no directament objecte d'estudi en aquest document* - que permeten la consecució dels objectius del projecte.

El desenvolupament d'aquesta plataforma estarà basat en la distribució de *Cloudera Apache Hadoop*.

3.3 Orientació de Negoci

L'arquitectura que la companyia seleccioni definirà en gran mesura les capacitats d'anàlisi del seu negoci. Aquest tipus de projectes, com passa amb els relacionats amb la disciplina del *Business Intelligence*, han de ser considerats i gestionats com a part de l'estratègia de la companyia.

Considerar un projecte d'aquestes característiques exclusivament des de l'àmbit de les TI, condemnarà al fracàs la iniciativa i malbaratarà l'esforç que, des del vessant tecnològic, s'hagi pogut dur a terme en aquest sentit.

Si bé el que ens ocupa és un projecte de fi de carrera, on bona part del pes específic es centrarà en les condicionalitats tecnològiques en les que es subscriu la feina aportada. La meua experiència en el món laboral així com en l'adopció i implantacions de solucions orientades a la millora del rendiment econòmic de les empreses, ha condicionat la inclusió d'aquest capítol en aquest treball per deixar palesa la necessitat d'una alineació i complementarietat de responsabilitats entre el negoci i la tecnologia per tal de garantir un èxit a tots els nivells en aquests tipus de projectes.

3.4 Escenaris reals d'Implantació per a una solució "Big data"

Gran part del que fa *Hadoop* així com d'altres tecnologies i enfocaments de tipus "Big Data" és que permeten a les empreses trobar respostes a preguntes que ni tan sols sabia que es podia plantejar.

Això pot donar lloc a idees que condueixin a noves idees o ajudar a identificar formes de millorar l'eficiència operativa. Actualment hi ha un nombre de casos en el que l'ús de grans volums de dades és crític.

En aquest punt s'exposen algunes orientacions de negoci on l'aplicació de tecnologies "Big data" serà d'ús obligat per a garantir el seu èxit i rendibilitat.

Motors de Recomanacions

Empreses online i els seus venedors fan ús de *Hadoop* per comparar i recomanar serveis o productes basats en l'anàlisi del perfil d'usuari i dades de comportament.

LinkedIn utilitza aquest enfocament per alimentar la seva característica de "Gent que podries conèixer", mentre que Amazon l'utilitza per suggerir productes per a la seva compra als consumidors en línia.

Models de Risc

Les empreses financeres, bancs i altres utilitzen Hadoop i Data warehouse de nova generació per analitzar grans volums de dades transaccionals i determinar el risc i l'exposició dels actius financers; d'aquesta manera poden avançar-se a possibles escenaris "*what-if*" basats en el comportament del mercat simulat i anotar els clients potencials en funció d'aquesta informació.

Detecció de Fraud

És possible emprar tècniques de "Big data" per combinar la informació del comportament del client, la seva informació històrica i les seves transaccions per tal de detectar activitat fraudulenta. Companyies de targetes de crèdit, per exemple, utilitzen aquests mecanismes per identificar el comportament transaccional que indica una probabilitat alta de detectar quan una targeta ha estat robada.

Anàlisi de Campanyes de Màrqueting

Els departaments de màrqueting de tots els sectors fa temps que utilitzen la tecnologia per monitoritzar i determinar l'efectivitat de les campanyes de màrqueting. El paradigma "Big data" permetrà als equips de màrqueting incorporar majors volums de dades cada vegada d'un nivell de granularitat major (tals com les dades originades de clics en ofertes o registres detallats de trucades) amb el propòsit d'augmentar la precisió de les tècniques d'anàlisi.

Anàlisi de Satisfacció de Clients

Les empreses de consum utilitzen *Hadoop* i les tecnologies relacionades amb “Big data” per integrar les dades dels canals d’interacció amb el client prèviament aïllats, com ara centres de trucades, xats en línia, twitter, etc, per obtenir una visió completa de l'experiència del client. Aquest seguiment permet a les empreses entendre l'impacte d'un canal d'interacció amb el client sobre un altre per tal d'optimitzar tot el cicle de vida de l'experiència del client.

Monitorització de Xarxes

Hadoop i altres tecnologies “Big data” també s’han emprat per consumir, analitzar i visualitzar les dades obtingudes dels servidors, dispositius d'emmagatzematge i altres equips de TI amb el propòsit que els administradors de sistemes puguin supervisar l'activitat de la xarxa i diagnosticar colls d'ampolla i altres problemes.

Aquest tipus d'anàlisi també es pot aplicar a altres formes de xarxes, tals com les xarxes de transport urbà, per exemple, per tal de millorar l'eficiència del combustible.

4 Disseny del Sistema

Com ja s'ha argumentat en els capítols anteriors, el paradigma “Big data” és diferent a d'altres filosofies de gestió de les dades, en molts aspectes.

Aquest capítol servirà per plantejar de quina manera s'ha realitzat el disseny del laboratori basat en aquesta tecnologia i quins condicionants s'han hagut de tenir en compte alhora de ser rigorosos en la seva construcció.

4.1 Decisions de disseny

El primer dels punts que s'ha tingut en compte alhora de plantejar un disseny d'aquestes característiques és el de ser coherent amb les especificacions que ha d'acomplir una plataforma d'aquesta índole. Amb el propòsit de presentar els avantatges que aquesta tecnologia aporta respecte a la **gestió, tractament i objectius** de la seva implantació, a la següent taula s'han disposat les diferents característiques.

Àmbit de les Dades	Estructura	Volum	Descripció
Dades Mestres	Estructurades	Baix	Aquelles entitats de dades que són de valor estratègic per a l'organització. Normalment no volàtils, i no transaccionals.
Dades Transaccionals	Estructurades i Semi Estructurades	Mig - Alt	Les transaccions comercials que es capturen durant les operacions i processos de negoci.
Dades de Referència	Estructurades i Semi Estructurades	Baix - Mig	Dades administrades internament o externament obtingudes per donar suport a la capacitat d'una organització a l'hora de processar eficaçment les transaccions, la gestió de dades mestres, i proporcionar capacitats de suport de decisions.
Metadada	Estructurades	Baix	Definides com les "dades sobre les dades". S'utilitzen com una capa d'abstracció per les descripcions estandaritzades i operacions.
Dades Analítiques	Estructurades	Mig – Alt	Dades derivades de l'operació del negoci. S'utilitzen per satisfer les necessitats de presentació d'informes i anàlisi.
Documents / Contingut On-line	No Estructurades	Mig - Alt	Documents, imatges digitals, dades geo-espacials, i arxius multimèdia.

Big Data	Multi Estructurades	Alt	Grans conjunts de dades que són un repte per emmagatzemar, buscar, compartir, visualitzar i analitzar.
-----------------	---------------------	-----	---

Figura 11. Taula comparativa sobre les diferents fonts que trobem actualment.

Aquestes diferents característiques influeixen en la manera de capturar, emmagatzemar, processar, recuperar i protegir la informació i evolucionar cap a un format “Big Data” i garantir-nos aquestes capacitats.

Novament presentem les decisions de disseny **en comparació** amb les que empràriem per qualsevol altre disseny arquitectònic d’un sistema d’informació ja conegut.

Àmbit de les Dades	Seguretat	Emmagatzematge / Accés	Modelatge	Processament / Integració	Consum
Dades Mestres Dades Transaccionals Metadada Dades Analítiques	Base de dades, aplicacions, usuaris - rols	SGBDR / SQL	Model predefinit (Relacional o Dimensional)	ETL, Replicació	BI, Eines estadístiques, aplicacions operacionals
Dades de Referència	Seguretat basada en la plataforma	XML / xQuery	Flexible i Extensible	ETL	Consum de dades basat en eines propietàries
Documents i Contingut On-line	Basat en el sistema de fitxers	Sistema de Fitxers / Cerca	Formularis Lliures	A nivell de SO	CMS
Big Data -Weblogs -Sensors -Xarxes Socials	Basat en el sistema de fitxers / Base de dades	Sistema de Fitxers Distribuït / no SQL	Flexible (valor clau)	Hadoop, MapReduce, ETL	BI, Eines Estadístiques

Figura 12. Taula comparativa sobre les tecnologies associades a la gestió de les dades presentades a la taula anterior.

Així doncs, l’elecció d’aquesta proposta donat, l’àmbit de les dades que volem gestionar, enmarca les pautes que cal seguir per oferir una solució tecnològica al respecte.

L'estructura estàndard per a una plataforma de tipus "Big data" es compon per una pila de tipus SMAQ¹⁵ que comprèn inicialment un total de 3 capes que poden ser ampliades en funció de les nostres necessitats.

Els sistemes de tipus SMAQ són normalment de tipus *open source*, distribuïts, i poden executar-se sobre hardware no certificat. Per al nostre sistema descriurem i introduïrem tot seguit quins són els elements que la configuraran.

La següent figura presenta per una banda la visió funcional de la pila: és a dir, quins són els elements de l'arquitectura que necessitem per a gestionar cada part de la plataforma, i relacionar visualment, els elements tecnològics que s'han implantat per respondre a les necessitats funcionals.

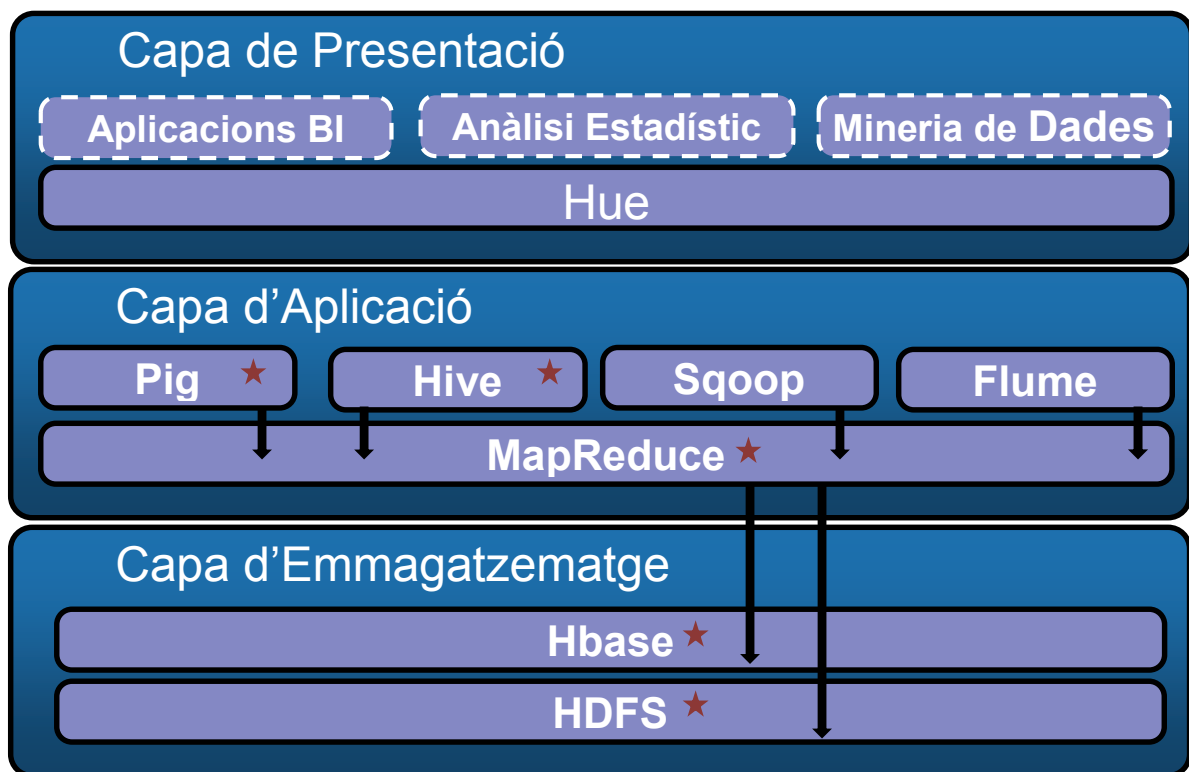


Figura 13. Pila "Big data" implementada per a la consecució dels objectius d'aquest projecte.

¹⁵ SMAQ; Es l'acrònim per: Emmagatzematge (Storage), MapReduce i Consulta (Query)

S'ha traslladat doncs a una arquitectura per capes el model “Big data”, de manera que ens serà més senzill identificar cadascuna de les funcionalitats amb la implementació corresponent. En el capítol 5, on s'amplien les característiques tècniques dels elements que conformen l'arquitectura, ens referirem a aquesta distribució per fer referència als diferents components de la plataforma.

Els elements ★ conformen la pila de tipus SMAQ.

Els elements — són els que comunament completaran una implantació d'aquests tipus de plataforma en un entorn productiu, si bé no formen part de l'abast d'aquest projecte.

4.2 Arquitectura

A l'apartat anterior s'han justificat una sèrie de decisions de disseny que afectaran la implementació de l'entorn “Big data” que es presentarà com a prova pràctica adjunt al present estudi.

En el punt que ens ocupa, cal detallar quines activitats realitzaran i de quina manera es comunicaran cadascun dels components que s'han escollit per tal de garantir la funcionalitat de la configuració proposada. La descripció dels elements que componen la pila “Big Data” descrita en la figura 13 es realitzaran mitjançant una descripció *Bottom – Up*.

En primer lloc es tractarà la capa física i els sistemes de d'emmagatzematge que propiciaran la base de la plataforma; tot seguit es detallaran els components que pertanyen a la capa d'aplicació i que, per aquests tipus d'arquitectures, tenen una importància vital a l'hora de gestionar les dades prèviament emmagatzemades: en aquest tipus d'arquitectura és possible controlar events que en una arquitectura habitual formen part exclusivament de la capa física; i finalment es farà una ullada al principal element de la capa de interfície d'usuari per tal de fer el seguiment i control dels processos a nivell administratiu.

4.2.1 Capa d'Emmagatzematge

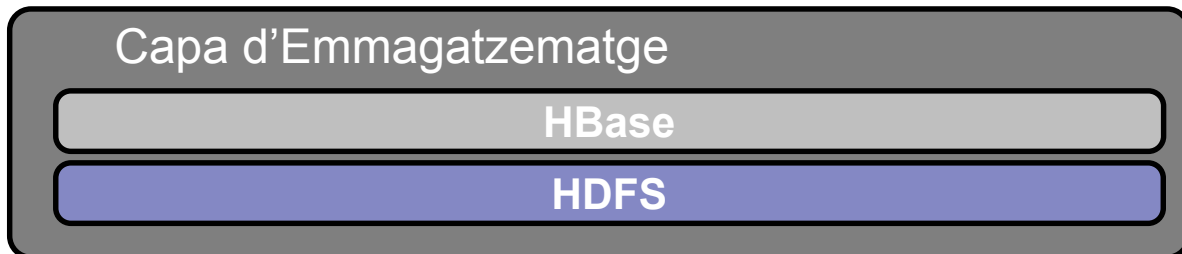
La disposició del nivell més baix de la solució determinarà el seu funcionament.

HDFS

Si es dona la situació en què un conjunt de dades creix més que la capacitat d'emmagatzematge d'una única màquina física, es fa necessari realitzar – d'alguna manera – una partició de la informació sobre un nombre de màquines separades.

A aquells sistemes de fitxers que són capaços de gestionar l'emmagatzematge d'informació a través d'una xarxa de màquines s'anomenen *sistemes d'arxius distribuïts*. HDFS està dissenyat per emmagatzemar arxius de grans dimensions mitjançant models **d'escriptura única** i lectura múltiple, que s'executen en clústers.

Els arxius suportats sota HDFS es divideixen en blocs de mida petita – la mida per defecte és de 64 MB – que s'emmagatzemen com a unitats de dades independents.



Un clúster HDFS té dos tipus de nodes que operen mitjançant un patró mestre-treballador:

Nodes Mestres

- El **NameNode**, gestiona l'espai de noms del sistema d'arxius. S'encarrega de mantenir l'arbre de fitxers així com les metadades de tots els arxius i directoris que conformen l'estructura de dades.
- El **JobTracker**, té com a funció llançar les tasques de MapReduce cap al nodes específics del clúster; és a dir, aquells que contenen les dades que cal tractar.

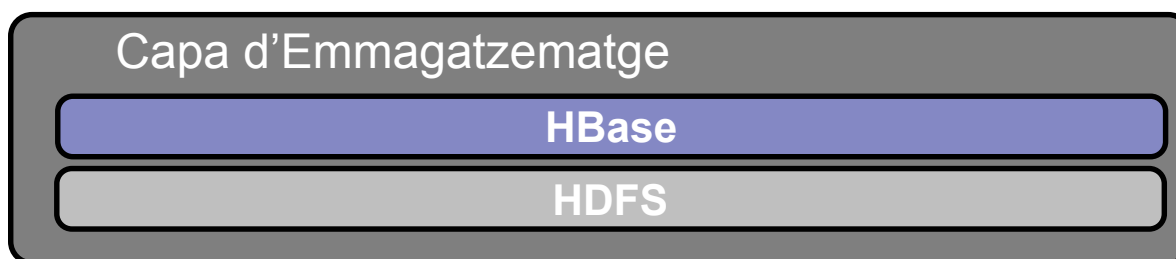
Nodes treballadors

- El **DataNode**, és el cavall de batalla del sistema d'arxius. Ells són els que s'encarreguen d'emmagatzemar i obtenir els blocs on s'ubiquen les dades – ja sigui per instrucció d'un NameNode o mitjançant una operació d'entrada-sortida generada pel client – i s'encarreguen de mantenir actualitzat al **NameNode** amb les llistes de blocs que estan emmagatzemant.
- Els **TaskTracker**, són un tipus de nodes preparats per acceptar tasques. Aquestes tasques són les operacions específicament executades per la funcionalitat MapReduce: Mapejar (Map), Reduir (Reduce) i Barrejar (Shuffle) – assignades per un node mestre de tipus **JobTracker**.

HDFS és el sistema d'emmagatzematge primari utilitzat per les aplicacions de tipus *Hadoop*. La característica principal d'aquest sistema de fitxers distribuït es la rapidesa en què és capaç de redirigir peticions a diferents nodes de computació a l'hora de realitzar càlculs sobre les dades emmagatzemades.

Hbase

El fet que HDFS sigui un sistema de fitxers orientat únicament a la d'addició, la necessitat de poder realitzar modificacions sobre aquelles dades prèviament emmagatzemades esdevé crítica per manipular la informació. HBase neix de la necessitat de disposar d'accés de lectura / escriptura en temps real sobre grans conjunts de dades.



HBase és una base de dades de tipus noSQL¹⁶, distribuïda – mantenint coherència amb l'arquitectura del sistema de fitxers que el suporta – i orientada a columnes¹⁷. La seva arquitectura està basada en BigTable¹⁸. HBase es caracteritza per les següents propietats:

- Afavoreix la **consistència** de dades vers l'**accessibilitat**.
- S'integra perfectament en qualsevol arquitectura basada en *Hadoop* (alta eficiència en càrregues massives i anàlisi MapReduce)
- Particions de Rang Ordenades (no fa ús d'algorismes Hash)
- Fragmenta i Escala la informació de forma automàtica.

HBase està modelat mitjançant l'ús d'un node mestre (*HBase*) que gestiona un grup d'un o més nodes esclaus (*RegionServer*). El node mestre és el responsable d'inicialitzar el sistema per a una nova instal·lació i d'assignar els espais de treball dels *RegionServers* registrats així com de la recuperació d'errors.

¹⁶ *noSQL: sistemes gestors de base de dades no relacionals*

¹⁷ *Paradigma d'emmagatzematge on – simplificant – s'optimitza l'accés a la informació indexant-la mitjançant columnes en comptes de fileres.*

¹⁸ [*Sistema de base de dades distribuïda*](#) creada per Google el 2004 i basada en GFS (Google File System)

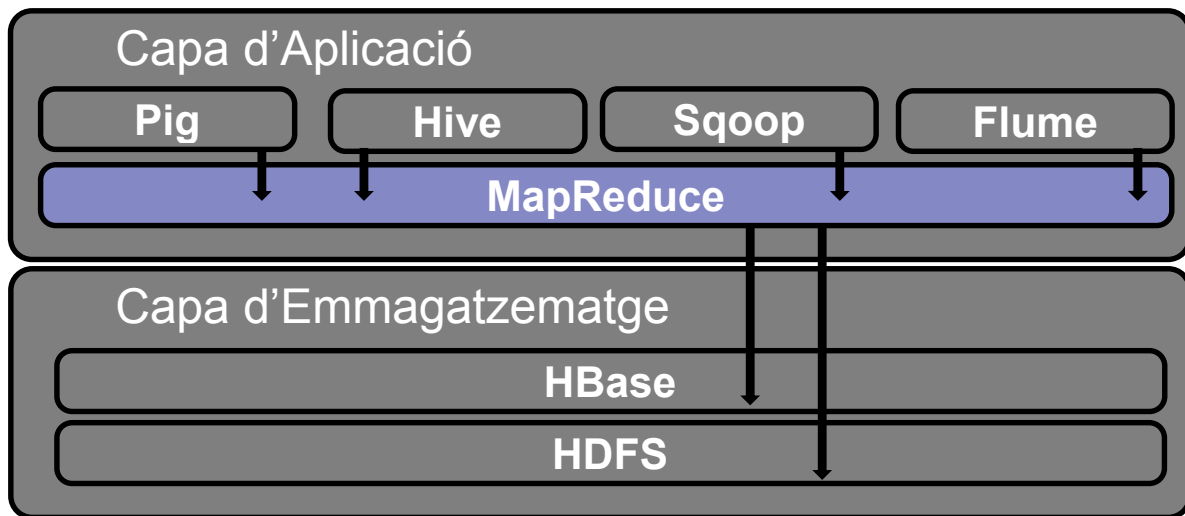
4.2.2 Capa d'Aplicació

La segona de les capes en què quedarà segmentat el nostre sistema és la que oferirà la capacitat de processament de la informació que prèviament ens hem encarregat d'emmagatzemar en el nostre sistema "Big data".

MapReduce

MapReduce és un *framework* que funciona com a capa de còmput per a informació emmagatzemada en fitxers distribuïts. Les tasques en les quals es separa la feina de MapReduce són quatre:

- La funció de mapejar (*Split and Map*), divideix les consultes sol·licitades pel client en diferents parts i s'encarrega de processar-les a nivell de node.
- La funció reduir (*Shuffle and Reduce*), agrega els resultats generats per la funció *Map* i determina la resposta final per a la consulta.



En particular trobem dos tipus de nodes que controlen aquest procés d'execució:

- Un *JobTracker* que s'encarrega de coordinar tots els treballs que s'executen en el sistema mitjançant la programació de tasques per executar-se sobre els *TaskTrackers*

- Un conjunt de *TaskTrackers*, que s'encarreguen d'executar les tasques sol·licitades pel *JobTracker* a qui remetent la informació generada i que li permet mantenir un registre de l'evolució general de cada treball. Si una tasca falla, el *JobTracker* pot reprogramar en un *TaskTracker* diferent per dur-la a terme.

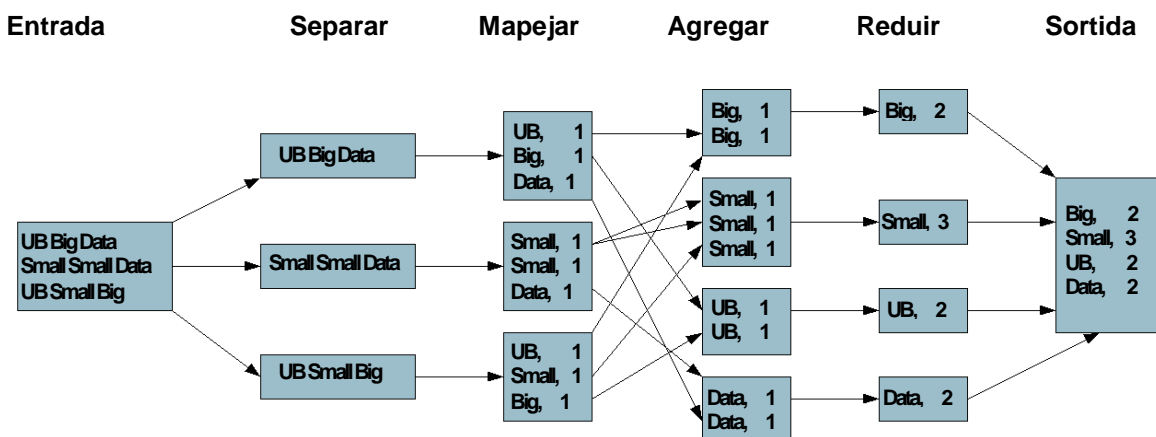
Aquest darrer punt, la seva **tolerància als errors**, és un dels avantatges més poderosos d'aquest sistema, que no només es capaç de processar la informació sino que també pot gestionar les excepcions automàticament per seguir refent els càlculs proposats per les peticions del client.

MapReduce. Un exemple pràctic.

De la correcta selecció i aplicació dels algorismes utilitzats en el nostre procés de *MapReduce* dependrà en gran part l'èxit en el processament de les dades implicades en la nostra solució "Big data".

Així doncs, si bé aquests algorismes poden variar segons les necessitats que tinguem a l'hora de programar els nostres processos *batch*, conceptualment és possible representar el procés a partir d'un input donat fins als resultats obtinguts.

Imaginem que tenim un conjunt de paraules separades per espais i volem analitzar el nivell d'aparició de cadascuna d'elles. Per al nostre process *batch* de MapReduce la part "Map" produirà parells del tipus (paraula, valor), mentre que la part "Reduce" afegeix 1 per cada paraula trobada, a més a més el nostre programa pre-reduceix en un pas posterior al de "Map" els resultats intermigs que finalment permet reduir el conjunt de sortida a un a més gran nivell.



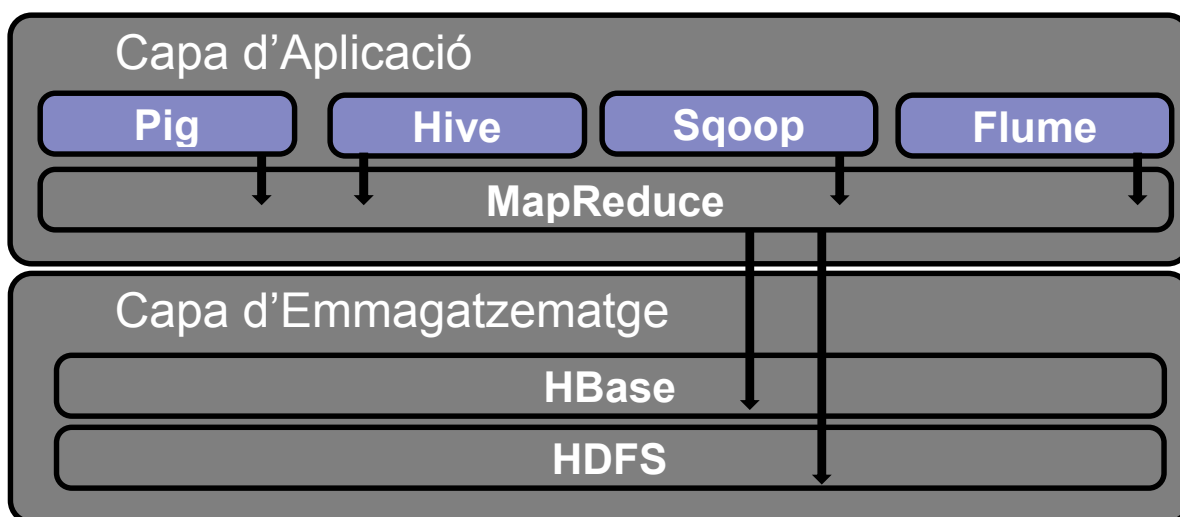
Procés 14. Procés *batch* MapReduce per al tractament d'un cas d'exemple.

Quan parlem de *Hadoop*, com a plataforma per a l'emmagatzematge i processament de dades, estem, en realitat, parlant de la funcionalitat complementària de dos components que ja hem presentat i conformen el *core* d'una plataforma "Big data":

HFDS + MapReduce = Hadoop

Hadoop és doncs capaç d'oferir la base tecnològica per al tractament dels casos que s'apliquen a models "Big data". Una altra característica que situa a *Hadoop* en una disposició prominent en aquest camp és que l'arquitectura és capaç d'**escalar linealment** ja sigui en funció de l'**increment de les dades** gestionades o bé a la **complexitat de les dades** processades.

Associat a aquesta capa trobem un grup de components que ofereixen capacitats funcionals als mecanismes base de l'arquitectura. Els presentarem i comentarem per tal de clarificar l'espai que ocupen en el conjunt de la solució.



Pig

Pig és un llenguatge de programació d'alt nivell per a generar programes *batch* de tipus MapReduce emprats per *Hadoop*.

El llenguatge emprat rep el nom de **Pig Latin** i, per aquells perfils familiaritzats amb llenguatge SQL té una gran similitud.

Com aquesta mena de llenguatge, Pig Latin pot ampliar-se fent ús de funcions (UDF¹⁹) que poden ser implementades en llenguatges procedurals o orientats a objectes tals com Java, Python o JavaScript i fer crides des d'aquest llenguatge. Els programes generats acabaran sent executats per les tasques MapReduce.

¹⁹ Acrònim en anglès per *User Defined Functions*

Hive

El segon dels components és fruit de la necessitat d'una empresa que necessitava poder tractar gran quantitat d'informació en un format analític similar al que pot proporcionar un data warehouse. El departament d'aquesta empresa va desenvolupar **Apache Hive** amb la idea de construir una infraestructura de tipus DWH que funcionés sobre *Hadoop* i obtenir la potència d'aquests tipus d'eines analítiques. Aquesta empresa, nascuda de l'emprenedoria d'un estudiant de Harvard el 2004, genera avui dia uns beneficis per valor de 3.71 bilions de dòlars ²⁰. Aquesta empresa és *Facebook*.

Aquesta capa d'abstracció permet realitzar consultes en un pseudo-llenguatge SQL anomenat HiveQL, que es converteix a *MapReduce* per ésser interpretat. Això permet a desenvolupadors sense coneixements sobre *MapReduce* fer ús del data warehouse i poder-lo integrar fàcilment amb diferents eines de Business Intelligence i de visualització de dades com *Microstrategy*, *Tableau*, etc.

Sqoop

Carregar grans quantitats de dades sobre *Hadoop* des dels diferents sistemes de producció pot esdevenir una tasca de difícil consecució. De quina manera podríem moure dades des de emmagatzematges aliens a la plataforma “Big data” i incorporar-los – o extreure'ls! – sobre el nostre sistema de fitxers (HDFS), Hive o Hbase? **Apache Sqoop** és la resposta.

Sqoop permet importar i exportar dades des de fonts estructurades tals com SGBDR, DWH i sistemes NoSQL. Els data-sets que s'incorporen s'escapen en particions i una tasca “MapReduce” de tipus “Map” mou aquesta informació i les ubica en els nodes corresponents.

²⁰ ["Facebook Current Report, Form 8-K, Filing Date July 26, 2012"](#). Font *SECDatabase.com*. 26 de Juliol de 2012.

Flume

El darrer *framework* que incorporem a la plataforma i que volem fer notar en aquest document es Flume.

Flume ens permetrà connectar Hadoop amb el món exterior i que qualsevol companyia té a la seva disposició en un entorn corporatiu. Així doncs serà possible incloure dades tals com aquelles que pertanyen a servidors web, aplicacions corporatives o dispositius mòbils per tal d'integrar-les a Hadoop.

4.2.3 Capa de Presentació

Finalment, i en el nivell d'abstracció més elevat tenim aquells elements que ens permetran interactuar amb la solució.

En el cas que ens ocupa, però, tindrem unes eines que, a diferència d'una solució estàndard, no estan orientades a l'usuari final. De fet l'element que descriurem és una eina de **monitorització de la plataforma**.



Hue

Hue ofereix una interfície web per a *Hadoop*, i pot ser utilitzat com una plataforma per a construir aplicacions basades en aquest sistema de “Big data”.

Tal i com s’ha comentat en el punt 4.1, els elements marcats amb línies discontinúes seran necessaris en un entorn de producció, si bé l’abast de la feina realitzada no cobreix la integració d’aquests tipus de components.

5 Implementació del Sistema

5.1 Requeriments principals

Aquest punt s'encarrega de disposar les condicions tècniques mínimes – tant a nivell de software com a nivell de hardware – necessaries per al funcionament de la plataforma de laboratori que es construirà com aplicació pràctica per al projecte.

Els elements ressaltats són els escollits per configurar la plataforma.

Software

En primer lloc revisarem les condicions de treball envers el software que emprarem per a la instal·lació i configuració de la plataforma. Al ser tots els components basats en codi obert, la seva implementació actual està preparada per funcionar en entorns Linux.

A continuació es descriuen quins S.O – en base a la seva distribució – són aptes per suportar els components necessaris per construir el nostre laboratori “Big data”.

Taula de Sistemes Operatius

Sistemes Operatius	Versió	Plataforma
Distribucions Red Hat		
Red Hat Enterprise Linux (RHEL)	5.7	64-bit
CentOS	6.2	64-bit, 32-bit
Oracle Linux Enterprise	5.6	64-bit
Distribucions SUSE		
SUSE Linux Enterprise Server (SLES)	11 (Service Pack 1 o posterior)	64-bit
Distribucions Ubuntu/Debian		
Ubuntu	Lucid (10.4) (LTS)	64-bit
	Precise (12.4) (LTS)	64-bit
Debian	Squeeze (6.03)	64-bit

Tot seguit es disposen els components que formaran part de la plataforma “Big data” i que oferiran les funcionalitats *anteriorment descrites*²¹

Taula de Components

Component	MySQL	SQLite	PostgreSQL	Oracle	Derby
Oozie	5.5	–	8.4	11gR2	Per Defecte
Flume	–	–	–	–	Per Defecte
Hue	5.5	Per Defecte	–	–	–
Hive	5.5	–	8.3	11gR2	Per Defecte
Sqoop	*			*	

* Sqoop pot treballar sobre MySQL 5.1 i HSQLDB 1.8 o superior.

A més la plataforma requerirà disposar d’una versió del software *Java Development Kit (JDK)* igual o superior a la 1.6.0_8. En particular l’entorn de laboratori farà us de la versió 1.6.0_31.

Hardware

Les necessitats, quant al suport físic de la plataforma, poden variar en funció de la potència que el client desitgi obtenir dels seus processos i de la dimensió de cadascuna de les tres “V” que vulgui tractar.

A continuació, com a part de la investigació duta a terme dins de l’abast d’aquest projecte i seguint amb la coherència dels plantejaments exposats, orientats a oferir informació per aplicar aquesta tecnologia en l’àmbit empresarial, es proposen 4 “recomanacions base” a tenir en compte alhora de dimensionar la plataforma en funció del volum de la informació a tractar i de la càrrega de treball.

- Configuració D. (1U²²): Dos unitats CPUs de quatre nuclis, 8GB RAM, 4 unitats de disc (1TB o 2TB).
- Configuració C. (1U): Dos unitats CPUs de quatre nuclis, de 16 a 24GB RAM, 4 unitats de disc (1TB o 2TB).
- Configuració B. (2U): Dos unitats CPUs de quatre nuclis, de 16 a 24GB RAM, 12 unitats de disc (1TB or 2TB).

²¹ Capítol 4.2 Arquitectura.

²² 1U-7U. Es la notació estàndard per descriure els Racks en que es construirà un sistema, sent 1U el tamany més reduït i 7U el més gran. [Font Viquipèdia 2012](#)

- Configuració A. (2U): Dos unitats CPUs de quatre nuclis, de 48 a 72GB RAM, 8 unitats de disc (1TB or 2TB).

L'exposada és la notació estàndard i ha de servir de referència com a base per a implantacions en entorns productius.

L'entorn de laboratori amb les següents característiques:

Sistema Operatiu (Host)	Versió	Processador	RAM	DISC
Windows 7	Ultimate	Intel Core i7-920XM (2GHz, 8MB Cache)	16GB DDR3 RAM (8GB + 8GB)	237GB
Sistema Operatiu (Laboratori)	Versió	Processador	RAM	DISC
CentOS	6.2	Intel Core i7-920XM (2GHz, 8MB Cache)	8GB DDR3 RAM	20GB

L'entorn de laboratori s'ha generat en un equip virtualitzat basat en *VirtualBox 4.2.6* juntament amb *VirtualBox 4.2.6 Oracle VM VirtualBox Extension Pack*.

5.2 Construcció del laboratori de proves

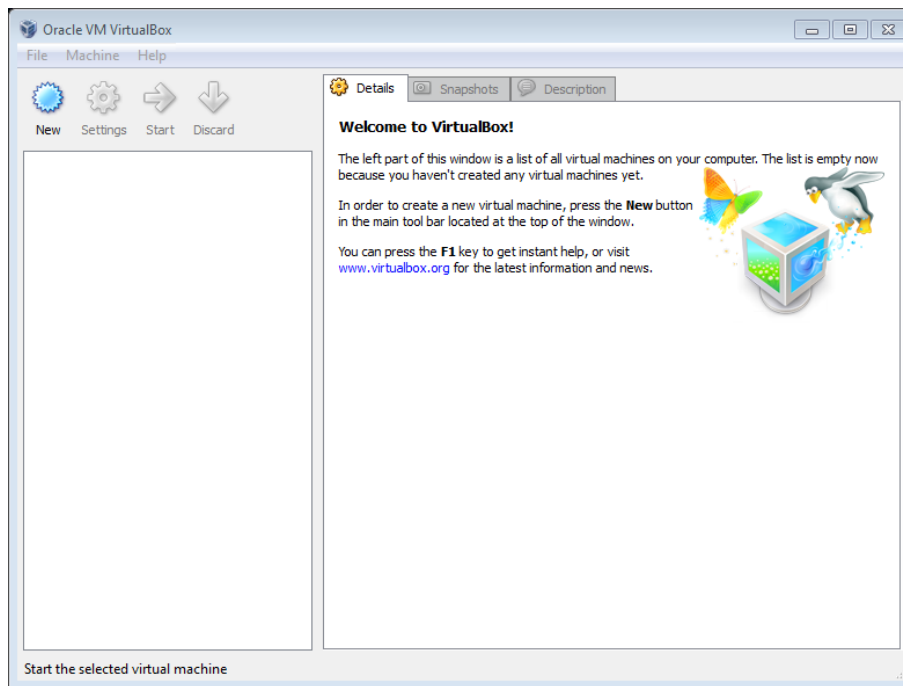
En aquest punt es descriu tot el procés de instal·lació i configuració de l'entorn "Big data" que s'ha plantejat com a objectiu per aquest Projecte Final de Carrera.

Cada pas s'etiquetarà en un sub capítol d'aquest apartat.

Per alleugerir el document tècnic, aquest estudi obviarà aquells passos que per ser repetitius i/o de notòria simplicitat seran eliminats fent referència a aquesta nota.

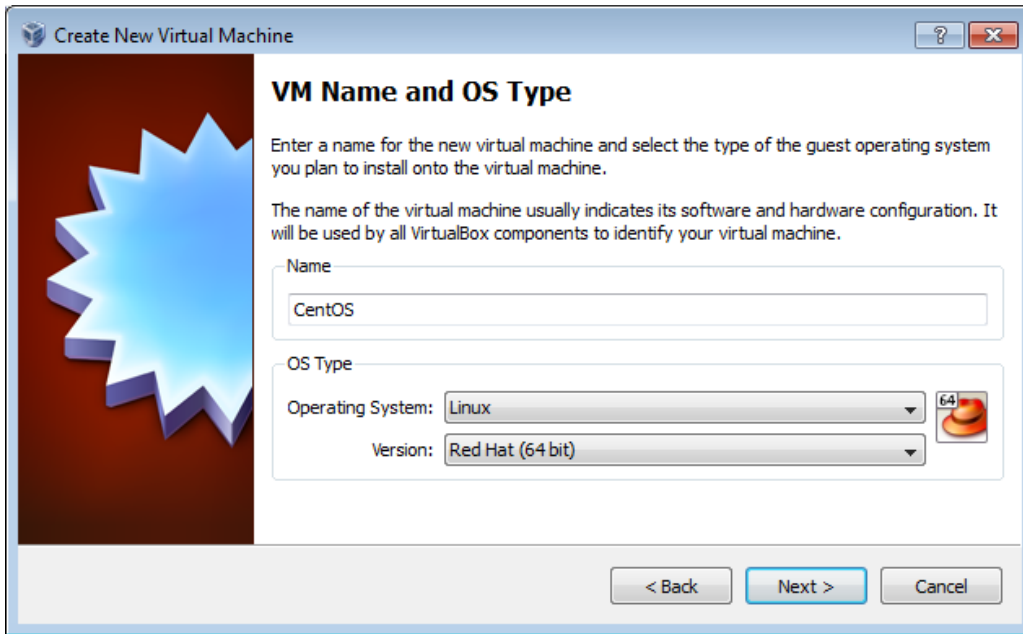
5.2.1 Instal·lació de CentOS 6.2 com a Sistema Operatiu Virtual

El primer pas que necessitem realitzar és instal·lar el nostre sistema operatiu en el seu format virtualitzat. Podem descarregar²³ la imatge *iso* del repositori oficial.

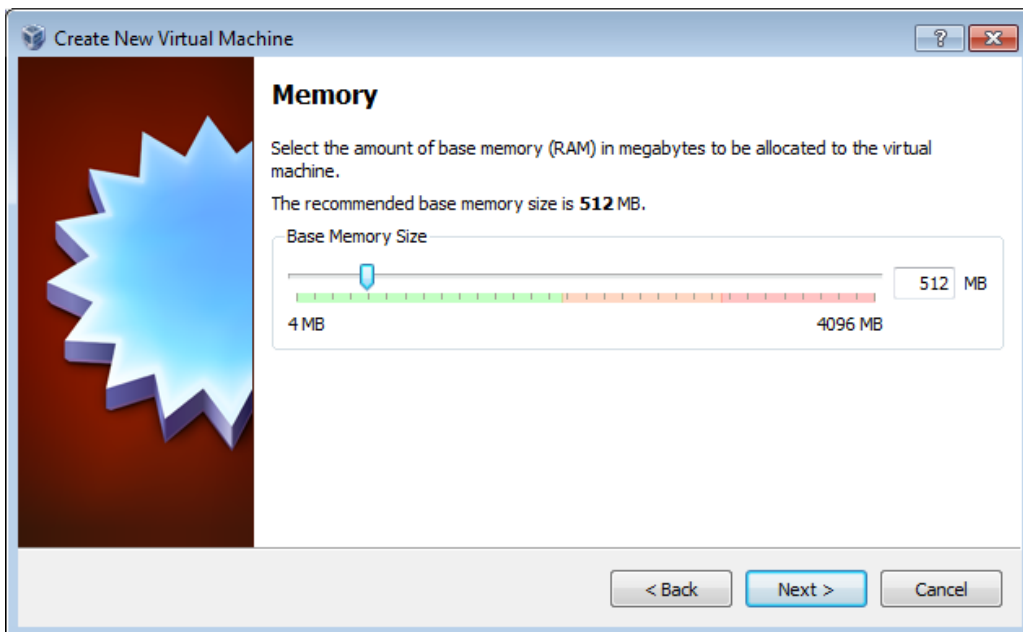


Escollim una etiqueta per a la nostra màquina virtual així com el sistema operatiu que instal·larem a continuació.

²³ Des del repositori oficial de CentOS es possible descarregar la darrera versió del sistema operatiu.
http://isoredirect.centos.org/centos/6/isos/x86_64/

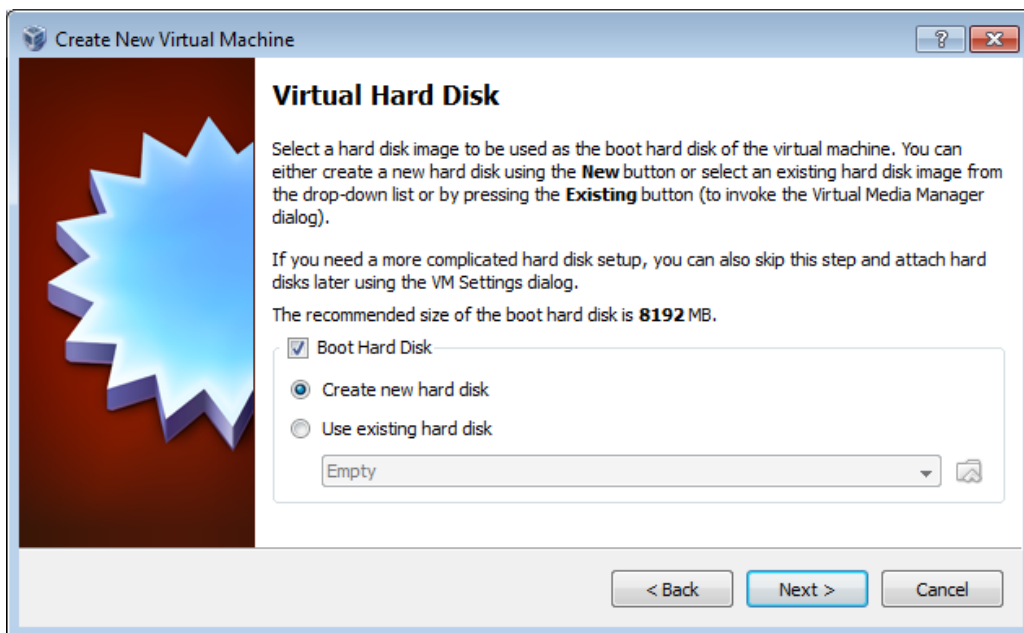


Triem la quantitat de memòria que desitgem assignar a la màquina virtual.

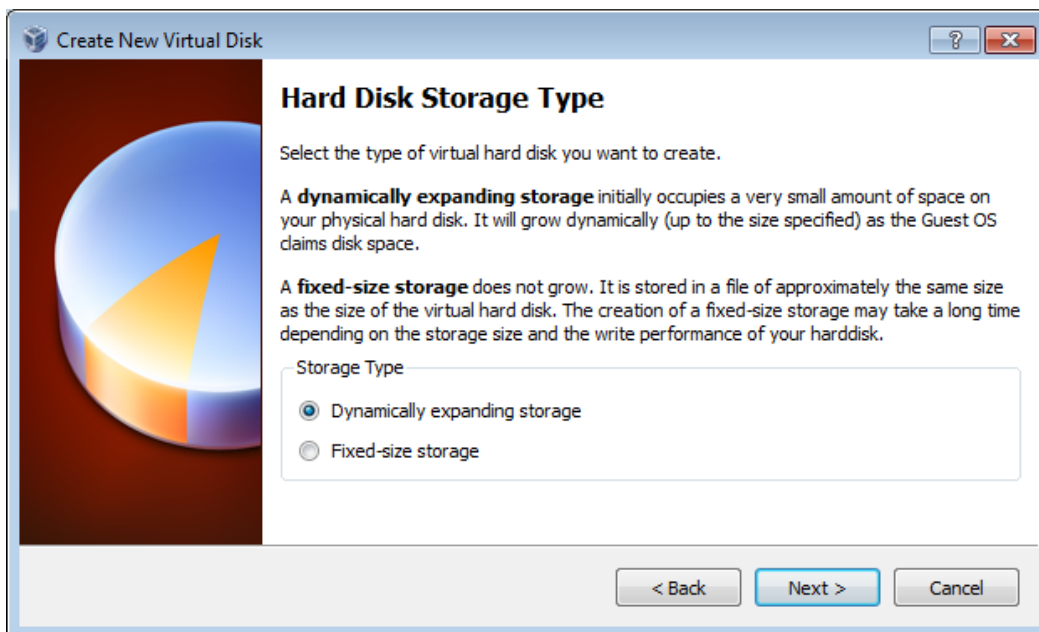


El següent pas serà triar el tamany del disc que contindrà tot el entorn.

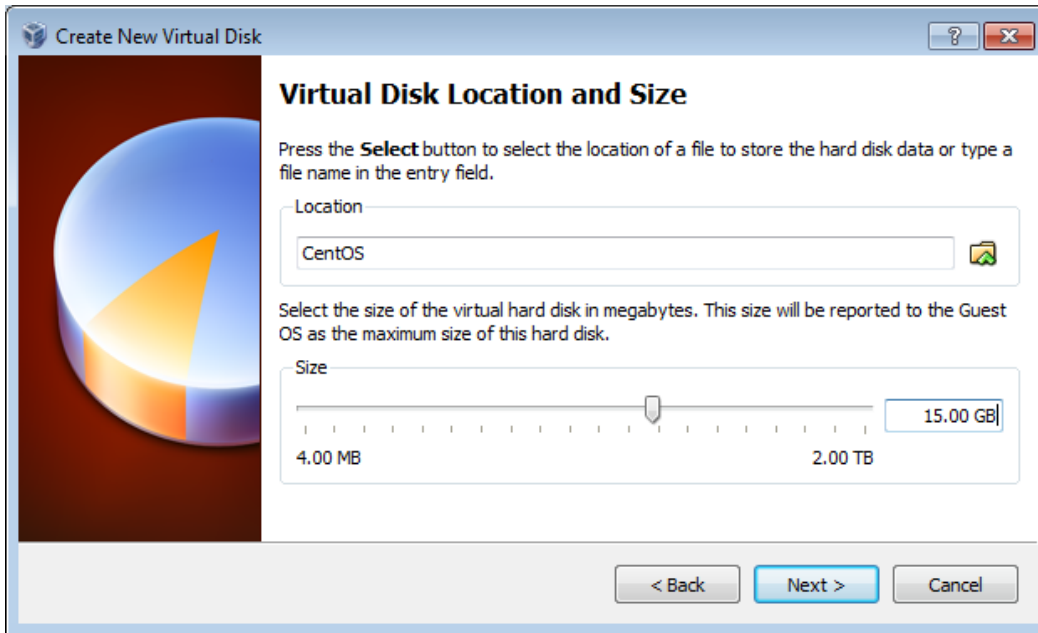
En aquest punt caldrà considerar l'espai ocupat pel sistema operatiu un cop instal·lat, i aquell que serà necessari per incorporar-hi tots els mòduls que conformaran la plataforma.



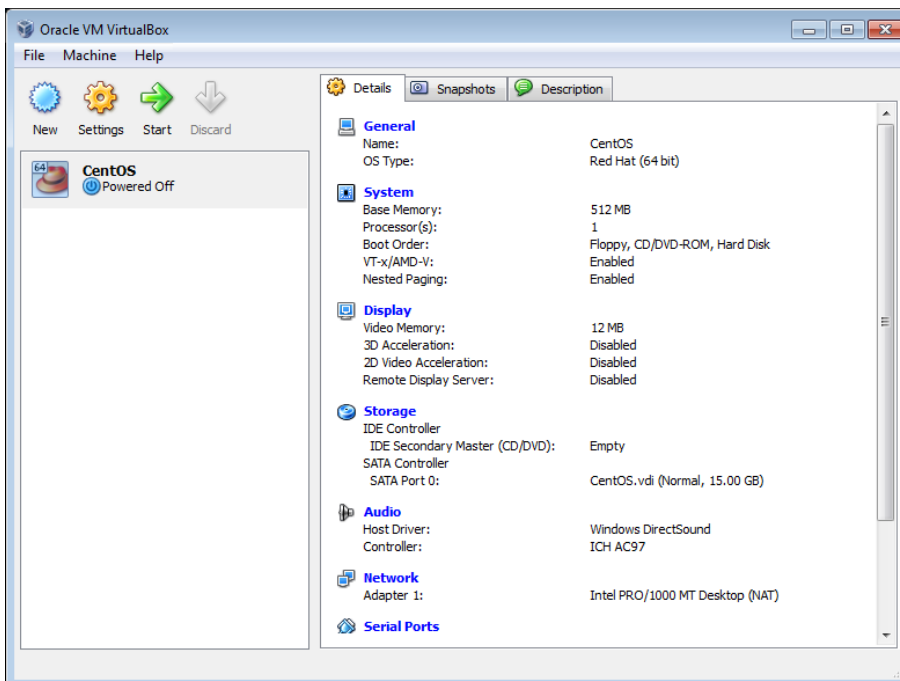
Escollim el tipus de creixement de la nostra màquina virtual. Pel nostre cas aquesta opció no resultarà crítica donat que el contingut de l'entorn estarà en tot moment acotat. Podem escollir l'opció que desitgem.



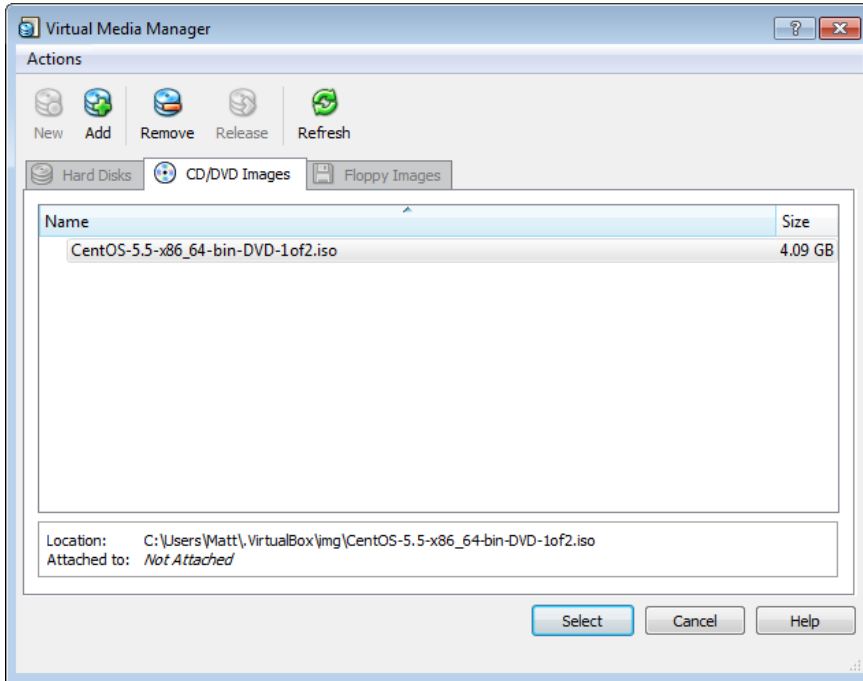
I en definirem el tamany inicial.



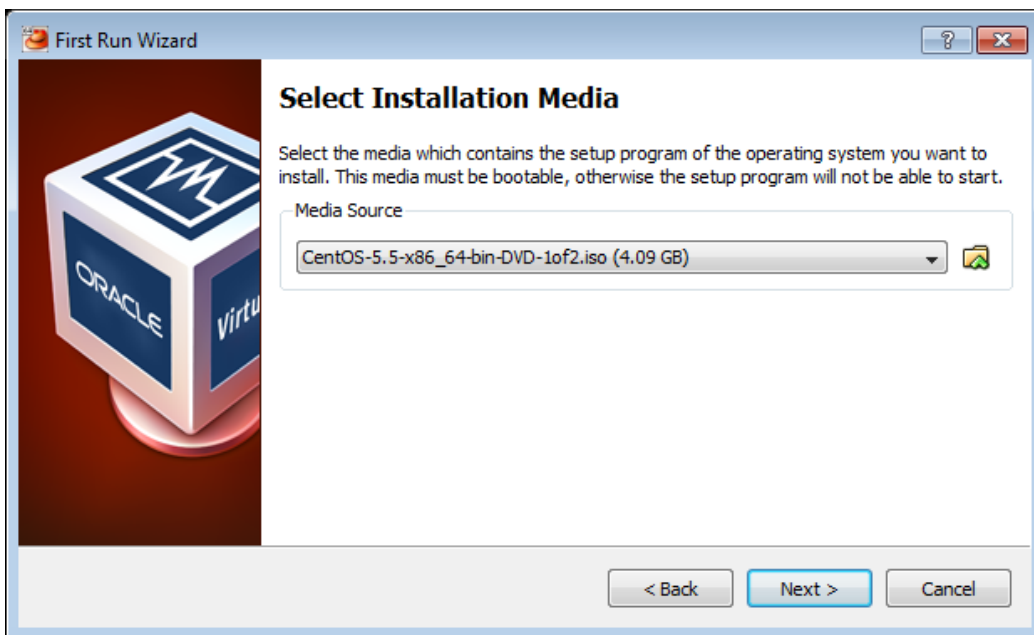
Finalment tindrem preparat l'entorn per a començar la instal·lació del sistema operatiu triat per al laboratori.



Per tal de fer-ho, escollirem la imatge que haguem descarregat i la configurarem com a base per a la màquina virtual, fet que propiciarà que el sistema operatiu arrenqui d'inici i permeti la seva instal·lació...



... i el seleccionem durant la primera arrencada de la plataforma.



El procés d'instal·lació del sistema operatiu serà el que ofereixi el programari per defecte mitjançant l'assistent. Per aquest motiu ometrem aquest pas.

5.2.2 Instal·lació del programari de suport

Tot seguit incorporarem aquell programari que serà necessari per garantir el funcionament de tot el sistema. Això és la instal·lació del kit de desenvolupament de Java de Sun²⁴ (Oracle JDK).

Instal·lació d'Oracle JDK

1. En primer lloc descarregarem la versió d'Oracle JDK des d'aquesta [adreça](#) fent clic a **Previous Releases** i a continuació a l'enllaç **Java SE 6**.
2. Com a usuari **root** ubiquem la variable d'entorn `JAVA_HOME` al directori on s'ha instal·lat el JDK; per exemple:

```
# export JAVA_HOME=<jdk-install-dir>
# export PATH=$JAVA_HOME/bin:$PATH
```

on `<jdk-install-dir>` seria quelcom similar a `/usr/java/jdk1.6.0_31`, depenen de la configuració del sistema on estigui instal·lat el JDK.

²⁴ Sun Microsystems fou adquirida per Oracle Inc el 20 d'Abril de 2009

5.2.3 Instal·lació dels components de la plataforma “Big data”

Els components que conformaran la plataforma són els descrits en el capítol [3.2 Comparativa i selecció d'eines](#) i dels que s'ha ampliat la informació en el capítol [4 Disseny del Sistema](#).

Són els següents: *Hadoop* (HBase + MapReduce), *Flume*, *Sqoop*, *Pig*, *Hive* i *Hue*

Instal·lació de Hadoop

Tal i com hem comentat en el capítol 3.2 aquesta plataforma es basa en l'arquitectura “Big data” proposada per un dels fabricants que actualment tenen una presència més notòria en aquest àmbit i que s'ha pres de referència per a l'elaboració d'aquest entorn.

En primer lloc, i en estar treballant sobre una distribució Red Hat, caldrà que afegim el repositori *yum* per instal·lar el CDH4²⁵ sobre el nostre sistema:

1. Afegirem el repositori *yum* a la nostra elecció, en el nostre cas:
 - a. `http://archive.cloudera.com/cdh4/redhat/6/i386/cdh/cloudera-cdh4.repo`
2. Descarregarem el fitxer: `cloudera-cdh4.repo` i el disposarem a la ubicació `/etc/yum/repos.d/`
3. Actualitzarem el paquet de *yum* executant
 - a. `$sudo yum update yum`
4. I cercarem el primer dels paquets que necessitem instal·lar: Hadoop
 - a. `$sudo yum search hadoop`
 - b. `$ yum install hadoop-0.20 -y`

Tot seguit, i amb el propòsit de simular una instal·lació en un entorn productiu, volem identificar amb una clau *hash* el nostre sistema dins un futurible cluster per tal de garantir que les comunicacions entre els nodes del nostre sistema d'informació seran segures.

Per fer-ho cal que afegim la clau corresponent des del repositori de Cloudera (Cloudera Public GPG Key) al nostre repositori amb la següent comanda:

²⁵ Acrònim per a *Cloudera Distribution Hadoop*


```
$ sudo rpm --import http://archive.cloudera.com/cdh4/redhat/6/x86_64/cdh/RPM-GPG-KEY-cloudera
```

Aquesta clau, a més, ens garantirà sempre l'accés a paquets originals d'actualitzacions per al nostre sistema.

La segona part de la instal·lació de *Hadoop* contempla la disposició de MapReduce.

Abans de continuar però, caldrà decidir on desplegar els *NameNodes*, i els *JobTrackers*.

- El **NameNode** i **JobTracker** s'acostumen a disposar – per regla general – en el mateix host “mestre” llevat que el cluster sigui més gran (més d'unes poques desenes de nodes), i el nostre host “mestre” no ha de executar el NameNode secundari, el DataNode o serveis TaskTracker.

Al ser el nostre laboratori una instal·lació “*stand alone*” aquesta primera consideració no és necessari tenir-la en compte.

- En clusters més grans, seria especialment important que el NameNode secundari s'executés en una màquina diferent a la del NameNode.
- Finalment, per aquests casos, cada node al cluster – amb excepció del host mestre – caldrà que executi el seu respectiu *DataNode* i els serveis *TaskTracker*

On cal fer la instal·lació?	Amb quines comandes?
JobTracker	<code>\$sudo yum install hadoop-0.20-mapreduce-jobtracker</code>
NameNode Primari	<code>\$sudo yum install hadoop-hdfs-namenode</code>
NameNode Secundari	<code>\$sudo yum install hadoop-hdfs-secondarynamenode</code>

Per a tots els clusters amb excepció

del JobTracker, NameNode

`$sudo yum install hadoop-0.20-mapreduce-tasktracker`

Primari i Namenode Secundari *

`hadoop-hdfs-datanode`

* No aplica per a un sistema *Stand-Alone*

En aquest punt ja tenim configurada la nostra plataforma Hadoop (Hbase+MapReduce) al sistema.

Instal·lació de Flume

En versions anteriors d'aquest component era necessari considerar certes dependències amb altres components de la plataforma.

En primer lloc, garantirem que no tenim cap instància del component instal·lada (òbviament en una instal·lació de zero aquest pas es pot ometre)

1. Garantirem que els processos que puguin estar funcionant sobre el servei – si existís – estan aturats abans de la nova instal·lació:
 - a. `$ sudo service flume- node stop`
2. Aturarem el servei en sí
 - a. `$ sudo service flume-master stop`
3. Desinstal·lem el paquet amb l'ús de *yum* executant
 - a. `$sudo yum remove flume`
4. I tot seguit instal·lem la darrera versió
 - a. `$sudo yum install flume-ng`

A més, desitjarem que el component estigui inclòs com a servei cada vegada que accedim a la plataforma.

5. Per tal de fer això caldrà instal·lar l'agent de Flume
 - a. `$sudo yum install flume-ng-agent`
 - b. `$sudo yum install flume-ng-doc`

Un cop instal·lat el component, la configuració es realitzarà de la manera següent.

Flume disposa d'una plantilla de configuració per al fitxer original *flume.conf* anomenat *conf/flume-conf.properties.template* i una altre plantilla per la fitxer *flume-env.sh* anomenat *conf/flume-env.sh.template*.

6. Realitzarem les còpies
 - a. `$ sudo cp conf/flume-conf.properties.template conf/flume.conf`
 - b. `$ sudo cp conf/flume-env.sh.template conf/flume-env.sh`

I parametritzarem els elements oportuns (canals de monitorització (TCP, UDP) , freqüència, etc.)

7. El següent pas es verificar la instal·lació,
 - a. `$ flume-ng help`
8. I finalment podem executar el servei (en les seves tres opcions)
 - a. `$ sudo service flume-ng-agent <start | stop | restart>`

Instal·lació de Sqoop

La instal·lació de Sqoop és aconsellable realitzar-la des del .tarball base del component, ja que els paquets d'instal·lació tenen en compte alguns factors que poden generar conflictes a l'hora de posar-ho en funcionament i ens obligaria a considerar certes dependències amb altres components de la plataforma.

1. La instal·lació és senzilla
 - a. `$ (cd /usr/local/ && sudo tar -zxvf _<path_to_sqoop.tar.gz>_)`
2. Tot seguit podem executar
 - a. `$ sqoop help`
 - b. `$ sqoop version`
 - c. `$ sqoop import`

Per tal de garantir la seva correcta instal·lació

3. Cal fer notar que en instal·lar Sqoop del paquet .tarball, cal garantir que les variables d'entorn `JAVA_HOME` i `HADOOP_MAPRED_HOME` estan configurades correctament.

La variable `HADOOP_MAPRED_HOME` ha d'apuntar al directori arrel de la instal·lació de Hadoop.

Opcionalment, i donat que també farem ús de les funcionalitats de HBase i Hive, també haurem d'assegurar que estan instal·lades les variables `HIVE_HOME` i `HBASE_HOME` perquè apunti al directori arrel de la instal·lació.

Instal·lació de Hue

Hue (*Hadoop User Experience*) és una interfície d'usuari Web per *Hadoop* que ofereix un conjunt d'aplicacions web i una plataforma per crear aplicacions personalitzades.

Per al funcionament de Hue, serà necessari instal·lar el paquet *hue-common* a la màquina on s'executarà. Si a més volem fer ús de MapReduce des d'aquest entorn caldrà instal·lar el paquet *hue-plugins* a les màquines on s'executin els *JobTrackers* (en el cas que ens ocupa serà la mateixa màquina)

1. En l'equip del servidor de Hue, instal·larem el meta-paquet de dades i el hue-server
 - a. `$ sudo yum install hue hue-server`
2. Tot seguit el paquet hue-plugins
 - a. `$ sudo yum install hue-plugins`
3. Configurar hue amb MapReduce
 - a. `cd /usr/share/hue`
 - b. `$ cp desktop/libs/hadoop/java-lib/hue-plugins-*.jar /usr/lib/hadoop-0.20-mapreduce/lib`
4. Iniciar hue
 - a. `$ sudo service hue start`

Amb això tindriem realitzada la instal·lació. Tot seguit es mostren les dependències respecte altres components que també formaran part de la plataforma:

Component	Requeriment	Aplicació
HDFS	Si	Core, FileBrowser
MapReduce	No	JobBrowser, JobDesigner, Beeswax
Hive	Si	Beeswax
HBase	No	Shell
Pig	No	Shell

S'han obviat configuracions avançades basades en autenticació segura, ldap, etc. que no s'han considerat en aquesta instal·lació per no ser objecte intrínsec del projecte.

Instal·lació de Pig

Seguint amb les clàusules comunes per a la instal·lació caldrà executar la següent sentència.

1. `$ sudo yum install pig`

Com en la instal·lació de Sqoop caldrà tenir en compte si volem fer visible aquest servei per a la resta de components (veure *Instal·lació de Sqoop*)

Amb això disposaríem de la plataforma consolidada.

6 Estudi de Costos i Viabilitat

Aquest punt conté, des de l'enfocament empresarial, una aproximació sobre el cost que podria representar per a una iniciativa privada l'esforç d'investigació i consolidació de la feina aquí realitzada per a l'adopció de la implantació descrita en aquest projecte.

6.1 Anàlisi del temps de realització del projecte

Per a la consecució dels objectius d'aquest projecte s'han emprat la major part de les hores de l'assignatura per a un total de 270 hores de documentació, donat que la que s'ha estudiat és una àrea relativament innovadora quant a l'aplicació empresarial es refereix i no ha estat senzill trobar referències per a casos d'ús reals.

La resta d'esforç 180 hores, han estat dedicades a la selecció dels components a consolidar i a preparar el disseny i implementació de la plataforma.

6.2 Valoració del cost econòmic del projecte

Es pretén calcular doncs el **cost d'exploració** amb l'objectiu de conèixer el cost de consecució del projecte i el cost unitari del producte (en el cas d'un cluster de treball es definirà un ràtio de comparació, per exemple "preu / hora de funcionament" o despeses de I+D+r d'un projecte de nou desenvolupament de plataforma, etc.)

El cost d'exploració es definirà com:

$$\text{CE (Cost exploració)} = \text{costos de producció} + \text{despeses generals.}$$

Per al càlcul dels costos d'exploració del sistema objecte d'aquest projecte, s'hauran de comptabilitzar les següents partides:

TAULA DE COSTOS D'EXPLORACIÓ	
Energía eléctrica	30 Kw. (0,115 € / kWh) x 8 hx 365 dies / anys = 10.074 € / any
Disseny Tècnic de la Solució	140 hores (100 €/hora) = 14.000 €
Implementació	180 hores (100€/hora) = 18.000 €
Contracte de manteniment i assistència tècnica	2.500 € /any
COST TOTAL (1er Any)	44.574 €
COST RECURRENT	12.574 € / Any

El cost total de l'exploració CE del sistema serà de : 44.574 l'any en que s'implanti la solució i de 12.574 € per als anys posteriors.

6.3 Retorn d'Inversió

Costa actualment dimensionar el ROI ²⁶ que genera l'incorporació d'aquestes pràctiques en la estratègia empresarial. Aquest fet és precisament el que fa que el guany no sigui prou palpable per a les direccions de les empreses i que la seva adopció quedi – a hores d'ara – restringida a corporacions amb un pressupost d'inversió molt elevat.

Les fonts consultades estimen el guany en un 35% en optimització de consulta de la informació o en l'increment de la capacitat que els centres de TI consoliden per a facilitar el consum i la claredat de la informació.

Per consultar alguns casos d'èxit en l'empresa privada es recomana llegir [aquest article](#) publicat per Rob Petersen amb data 17/12/2012.

²⁶ De l'anglès "Return Of Investment"

7 Conclusió

Per finalitzar, es revisa en aquest apartat l'objectiu del projecte, les fites assolides així com les principals línies de continuació del projecte en estadis futurs.

Objectius Globals

El projecte ha permès assolir els objectius proposats a l'inici. Això es:

- S'ha recopilat i contrastat informació al voltant de l'actual situació en la gestió de la informació empresarial (objectiu específic 1 i 2)
- S'han proposat les eines que millor cobreixen – en el moment de desenvolupar-se aquest Projecte Final de Carrera – les mancances que actualment es presenten en l'anàlisi i processament d'informació en temps real. (objectiu específic 3)
- S'ha dissenyat una arquitectura “Big data” basada en el software OpenSouce actualment disponible que s'integra amb l'espectre més gran de configuracions de sistemes que es troba a l'abast de les empreses de negocis (objectiu específic 4)
- S'ha construït un laboratori per mostrar la integració del sistema tal i com s'havia proposat inicialment. (objectiu específic 5)

Objectius Acadèmics

Acadèmicament la proposta m'ha permès recuperar coneixements assolits principalment en les assignatures de Sistemes Operatius (gestió de memòria - Jaume Timoneda), Minería de Dades (algorismes de reconeixement de patrons - Josep Fortiana) , Base de dades (Sistemes transaccionals, SQL - Anna Puig, Carles Franquesa) i Estructura de Dades.

Objectius Personals

Estic plenament convençut de l'encert en l'elecció de la temàtica del present Projecte de Final de Carrera que, de forma consensuada, el meu director – professor Enric Biosca Trias - i jo mateix vàrem plantejar.

Sóc de l'opinió que el tema d'anàlisi no pot ser més **actual**, tecnològicament **innovador** i que tot i trobar-se en un estadi inicial presenta signes de **maduresa** i és objecte ja d'**implantacions en el món real** per que és d'una importància capdal per aquells casos de negoci que la requereixen.

Tots aquests factors han estat alhora un *handicap* i una motivació afegida que m'han permès, no només consolidar el darrer pas per a l'obtenció de la meva titulació universitària, sinó generar un coneixement directament exportable al meu àmbit de treball i amb el qual serà possible enriquir tant el meu perfil professional com el port-foli de serveis i propostes tecnològiques que podran realitzar-se des del departament que actualment gestiono.

8 Referències Bibliogràfiques

Document	Editorial	Autor	ISBN
<i>Tapping into Unstructured Data: Integrating Unstructured Data and Textual Analytics into Business Intelligence</i>	<i>Prentice Hall</i>	<i>William H. Inmon (Author), Anthony Nesavich</i>	<i>0-13-236029-2</i>
<i>Hadoop: The Definitive Guide, 3rd Edition</i>	<i>O'Reilly Media, Inc.</i>	<i>Tom White</i>	<i>978-1-4493-1152-0</i>
<i>Big Data Analytics: Disruptive Technologies for Changing the Game</i>	<i>MC Press</i>	<i>Dr. Arvind Sathi</i>	<i>978-1583473801</i>
<i>Ethics of Big Data: Balancing Risk and Innovation</i>	<i>O'Reilly Media, Inc.</i>	<i>Kord Davis</i>	<i>978-1449311797</i>
<i>Big Data Analytics: Turning Big Data into Big Money (Wiley and SAS Business Series)</i>	<i>Wiley</i>	<i>Frank J. Ohlhorst</i>	<i>978-1118147597</i>
<i>Programming Hive</i>	<i>O'Reilly Media, Inc.</i>	<i>Edward Capriolo</i>	<i>978-1449319335</i>
<i>Hbase Definitive Guide</i>	<i>O'Reilly Media, Inc.</i>	<i>Lars George</i>	<i>978-1449396107</i>
<i>MapReduce Design Patterns: Building Effective Algorithms and Analytics for Hadoop and Other Systems</i>	<i>O'Reilly Media, Inc.</i>	<i>Donald Miner, Adam Shook</i>	<i>978-1449327170</i>
<i>NoSQL Handbook</i>	<i>http://nosqlhandbook.com/</i>	<i>Mathias Meyer</i>	

I. Annex. Manual d'Usuari de l'entorn de laboratori

Aquest apartat té com a propòsit presentar l'entorn de treball on l'usuari pot realitzar proves funcionals de l'arquitectura que s'ha implementat per a la consecució dels objectius proposats en aquest projecte.

L'entorn de treball presentat fa ús de les eines configurades per tal de simplificar la consulta d'informació referent a la plataforma.

Tota la informació que es mostra és susceptible de ser obtinguda mitjançant la consola de comandes, si bé l'alumne ha cregut convenient aprofitar els elements prèviament instal·lats i configurats per tal de presentar les eines disponibles.

A l'espai per a la disposició d'accessos directes es disposen 4 elements dels quals podem revisar-ne el seu funcionament:

Interfícies d'Usuari. Arquitectura

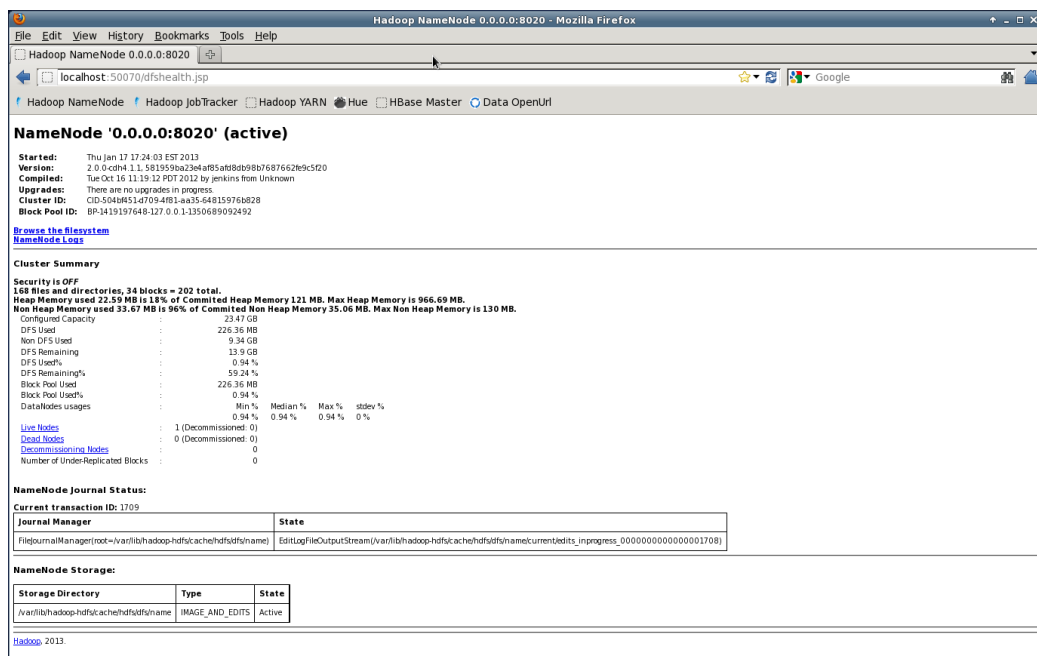
- *Hadoop NameNode*, pàgina jsp que consulta el sistema i retorna informació sobre els NameNodes configurats al sistema.
- *Hadoop JobTracker*, pàgina jsp que consulta el sistema per retornar informació sobre les tasques que s'han executat o s'estan executant en el sistema.
- *HBase Master*, HBase és el sistema de base de dades distribuït que suporta l'emmagatzematge estructurat per a grans taules. En l'entorn de laboratori està configurat en mode *standalone*. En aquest mode, la màquina virtual de java ofereix suport al HBase Master, i la capa de coordinació ZooKeeper.

Interfícies d'Usuari. Aplicació

- *Hadoop YARN*²⁷, interfície per a la gestió d'aplicacions MapReduce.
- *Hue*²⁸, és la interfície gràfica creada pel fabricant de referència en que es basa l'arquitectura Hadoop escollida (*Cloudera*) per mantenir un control integrat de tot l'ecosistema de la plataforma. Inclou (Beewax – Interfície per a consultes HiveQL, Gestor de Tasques i Oozie – el gestor de control de tasques –)

Hadoop NameNode

Realitzant una petició servidor web intern de Hadoop (port 50070) accedim al recurs `dfshealth.jsp`, una interfície web que ens permet consultar l'estat actual del sistema de fitxers distribuït.



The screenshot shows a web browser window displaying the Hadoop NameNode health page. The page title is "NameNode '0.0.0.0:8020' (active)". It provides detailed information about the NameNode's status, including start time, version, and configuration. A "Cluster Summary" section shows that security is off and provides metrics on heap memory usage and DFS capacity. Below this, there are sections for "NameNode Journal Status" and "NameNode Storage" with associated tables.

NameNode '0.0.0.0:8020' (active)

Started: Thu Jan 17 17:24:03 EST 2013
Version: 2.0.0-cdh4.1.1.581959ba294a4a85a88b098b76876629dc520
Compiled: Tue Oct 16 11:19:12 PDT 2012 by jenkins from Unknown
Upgrades: There are no upgrades in progress.
Cluster ID: CID:504885347094603-aa3a3c44815976a828
Block Pool ID: BP-3419197648-127.0.0.1-1350689092492

[Browse the filesystem](#)
[NameNode Logs](#)

Cluster Summary

Security is OFF
188 files and directories, 34 blocks = 202 total.
Heap Memory used 22.59 MB is 18% of Committed Heap Memory 121 MB. Max Heap Memory is 966.69 MB.
Non-Heap Memory used 33.67 MB is 96% of Committed Non-Heap Memory 35.06 MB. Max Non-Heap Memory is 130 MB.

Configured Capacity	23.47 GB
DFS Used	226.36 MB
Non-DFS Used	9.34 GB
DFS Remaining	13.9 GB
DFS Usage%	0.94 %
DFS Remaining%	59.24 %
Block Pool Used	226.36 MB
Block Pool Usage%	0.94 %
DataNodes usages	Min % Median % Max % stdev %
	0.94 % 0.94 % 0.94 % 0 %

[Live Nodes](#): 1 (Decommissioned: 0)
[Shard Nodes](#): 0 (Decommissioned: 0)
[Decommissioning Nodes](#): 0
Number of Under-Replicated Blocks: 0

NameNode Journal Status:
Current transaction ID: 1709

Journal Manager	State
filejournalmanager[not=/var/lib/hadoop-hdfs/achehdfsfsname]	EditLogOutputStream[=/var/lib/hadoop-hdfs/achehdfsfsname/current/edits_inprogress_000000000000001708]

NameNode Storage:

Storage Directory	Type	State
/var/lib/hadoop-hdfs/achehdfsfsname	IMAGE_AND_EDITS	Active

Hadoop, 2013

Figura 1. La interfície també permet navegar pel sistema de fitxers

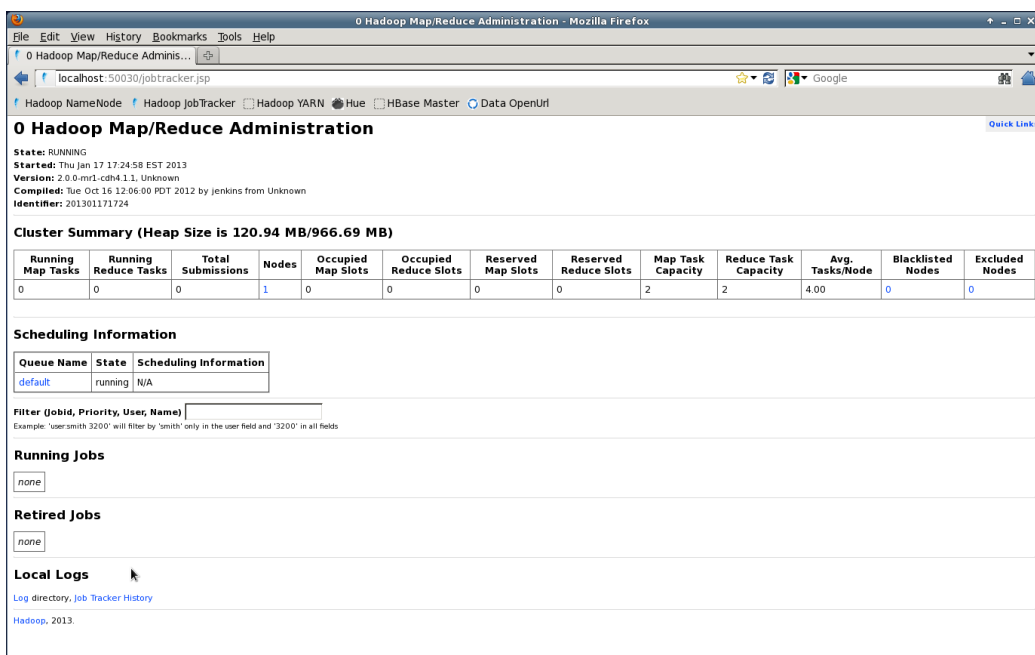
²⁷ YARN, es el nom que s'adopta per fer referencia a la versió 2.0 de MapReduce (MRv2)

²⁸ HUE, es l'acrònim per *Hadoop User Experience*.

Hadoop JobTracker

Cada Namenode i DataNode corren sobre un servidor web intern per tal de poder mostrar a l'administrador sobre l'estat actual del cluster. Amb la configuració que es presenta, la pàgina des de la qual podem consultar el NameNode es <http://namenode:50070/> i la llista del conjunt de DataNodes que conformen el cluster i les seves dades bàsiques.

El mateix servidor conté un altre recurs (port 50030) anomenat jobtracker.jsp que es el servei que s'encarrega de gestionar les operacions MapReduce executades en les adreces als nodes del cluster.



0 Hadoop Map/Reduce Administration

State: RUNNING
Started: Thu Jan 17 17:24:58 EST 2013
Version: 2.0.0-mr1-cdh4.1.1, Unknown
Compiled: Tue Oct 16 12:06:00 PDT 2012 by jenkins from Unknown
Identifier: 201301171724

Cluster Summary (Heap Size is 120.94 MB/966.69 MB)

Running Map Tasks	Running Reduce Tasks	Total Submissions	Nodes	Occupied Map Slots	Occupied Reduce Slots	Reserved Map Slots	Reserved Reduce Slots	Map Task Capacity	Reduce Task Capacity	Avg. Tasks/Node	Blacklisted Nodes	Excluded Nodes
0	0	0	1	0	0	0	0	2	2	4.00	0	0

Scheduling Information

Queue Name	State	Scheduling Information
default	running	N/A

Filter (Jobid, Priority, User, Name)
Example: "user:smith 3200" will filter by "smith" only in the user field and "3200" in all fields.

Running Jobs

Retired Jobs

Local Logs
[Log directory](#), [Job Tracker History](#)
Hadoop, 2013.

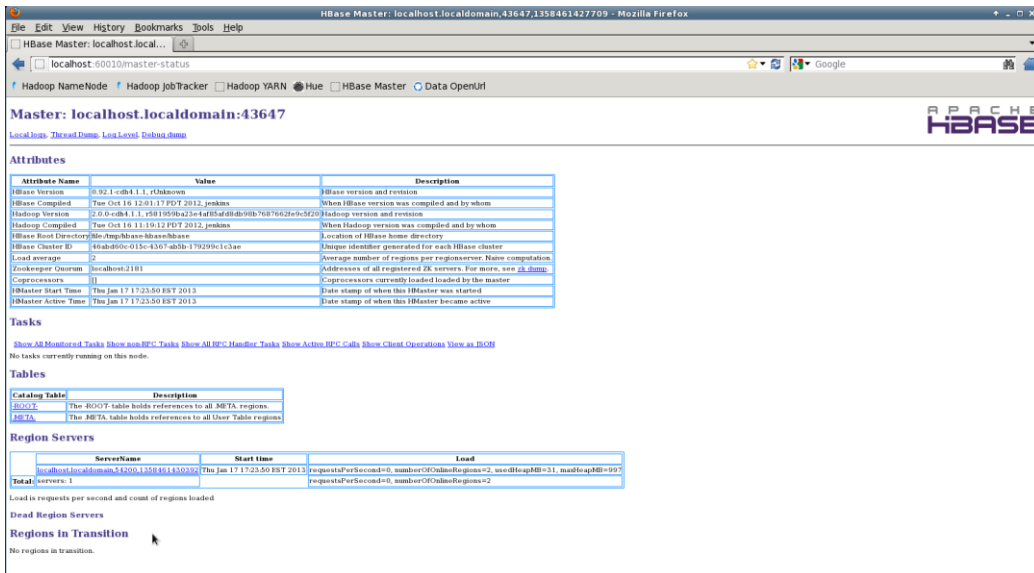
Figura 2. La interfície per la JobTracker permet controlar els paràmetres d'execució de les tasques MapReduce

Aquí podem observar informació del cluster així com el tipus de tasca en el seu moment d'execució juntament amb els recursos del sistema que consumeix.

Hbase Master (HMaster)

HMaster és la implementació del servidor mestre. El servidor mestre en un sistema emmagatzematge distribuït és el responsable de monitoritzar les instàncies dels seus RegionServers en el cluster, alhora que funciona com a capa de comunicació per tots els canvis que es realitzin sobre la metadada.

Per a controlar el seu funcionament tenim una interfície web (port 60010) que llista el conjunt de taules creades en el sistema així com la seva definició (ColumnFamilies, blocksize, etc.). De forma adicional – tot i que no en el cas que ens ocupa – si estan disponibles, també és possible navegar sobre els RegionServers del cluster i comprovar la seva configuració.



The screenshot shows the HBase Master web interface in a Mozilla Firefox browser. The page title is "Master: localhost.localdomain:43647". The interface includes a navigation bar with links for Hadoop NameNode, Hadoop JobTracker, Hadoop YARN, Hue, HBase Master, and Data OpenUrl. The main content area is divided into several sections:

- Attributes:** A table listing various system attributes and their values.
- Tasks:** A section indicating that no tasks are currently running.
- Tables:** A section listing catalog tables like 'hbase' and 'hbase:meta'.
- Region Servers:** A table showing the status of region servers, including server names, start times, and load metrics.
- Dead Region Servers:** A section indicating that no regions are in transition.

Attribute Name	Value	Description
HBase Version	0.92.1-cdh4.1.1 (Unknown)	HBase version and revision
HBase Compiled	Tue Oct 16 13:41:11 PDT 2012, Jenkins	When HBase version was compiled and by whom
Hadoop Version	2.0.6-cdh4.1.1, r5811059ba23e4af8af4d4db99b76876629e3c520	Hadoop version and revision
Hadoop Compiled	Tue Oct 16 11:39:12 PDT 2012, Jenkins	When Hadoop version was compiled and by whom
HBase Home Directory	/usr/lib/hbase-hbase	Location of HBase home directory
HBase Cluster ID	44ab49c-015c-4367-d3b-179299c1c3aa	Unique identifier generated for each HBase cluster
Load average	2	Average number of regions per regionserver. Name computation
Zookeeper Quorum	localhost:2181	Addresses of all registered ZK servers. For more, see zk.html .
Coprocessors	0	Coprocessors currently loaded by the master
HMaster Start Time	Tue Jan 17 17:23:50 EST 2013	Date stamp of when this HMaster was started
HMaster Active Time	Tue Jan 17 17:23:50 EST 2013	Date stamp of when this HMaster became active

ServerName	Start time	Load
localhost.localdomain:43647,1358481430350	Tue Jan 17 17:23:50 EST 2013	requestsPerSecond=0, numberOfOnlineRegions=2, usedHeapMB=31, modifiedMB=997
Total:	servers: 1	requestsPerSecond=0, numberOfOnlineRegions=2

Figura 3. Pàgina de consulta del HBase Master

YARN (MapReduce v2.0)

Com hem comentat en el detall de la implementació d'aquest projecte, Hadoop és el resultat d'una capa d'emmagatzematge (HDFS) i una capa de processament (MapReduce).

Aquesta capa de processament s'ha ampliat en versions posteriors per a dotarla de més capacitat adaptació a la varietat d'informació que Hadoop comença a tractar en diferents architectures. **YARN**, és el resultat de incloure en un únic framework les capacitats de MapReduce complimentades amb un gestor de recursos i aplicacions distribuït que permet generar processos a mida.

Així doncs MapReduce passa a ser des de la seva versió 2.0 un component de YARN.

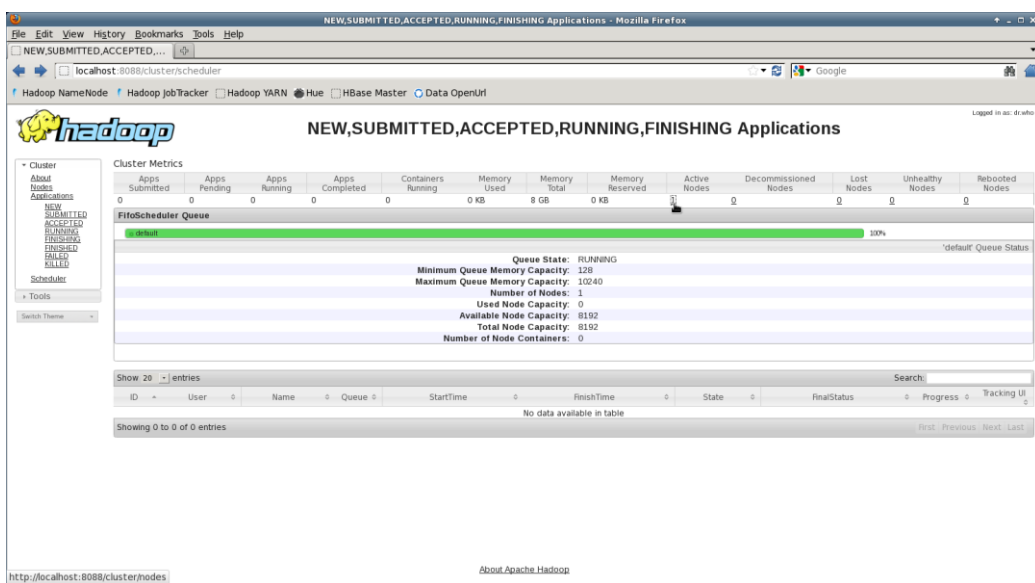


Figura 4. YARN permet controlar la execució d'aplicacions MapReduce.

Beewax

Beewax és la interfície gràfica per a l'execució de consultes HiveQL.

L'entorn permet parametritzar consultes, carregar informació i desar consultes ja executades.

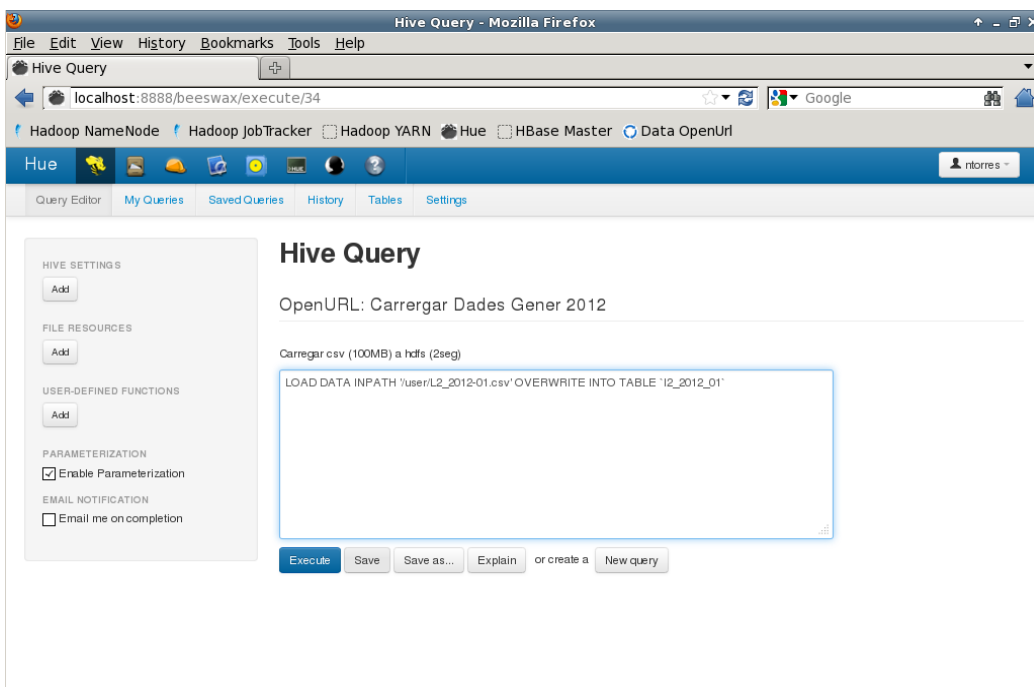


Figura 5. Una de les capacitats interessants de la interfície es la de generar automàticament les taules en les que s'emmagatzema la informació dels fitxers carregats.

Oozie

Finalment aquesta aplicació ens permetrà crear, modificar i programar fluxes d'execució de tasques MapReduce.

El quadre de comandament ens permetrà mantenir un seguiment en temps real de les tasques **agrupades** en fluxes de processos. És possible crear nous fluxes i programar-los convenientment per a la seva execució.

The screenshot displays the Oozie App interface in Mozilla Firefox. The browser's address bar shows 'localhost:8888/oozie/'. The interface includes a navigation bar with 'Hue', 'Hadoop NameNode', 'Hadoop JobTracker', 'Hadoop YARN', 'Hue', 'HBase Master', and 'Data OpenUrl'. Below this is a 'Dashboard' section with tabs for 'Workflows' and 'Coordinators'. A search filter is present, and a 'Show only' dropdown is set to '1' days with status options for 'Succeeded', 'Running', and 'Killed'. The 'Running' section shows a table with one entry: 'obterir_paraulas' with a status of 'RUNNING' and 50% progress. The 'Completed' section shows a table with six entries, including 'compla_paraulas' (KILLED), 'streaming_wordcount' (SUCCEEDED), and 'MapReduce' (SUCCEEDED).

Submission	Status	Name	Progress	Submitter	Id	Action
17 Jan 2013 01:34:09	RUNNING	obterir_paraulas	50%	rtorres	0000006-130116195359883-oozie-oozi-W	Kill

Completion	Status	Name	Duration	Submitter	Id
17 Jan 2013 01:20:35	KILLED	compla_paraulas	15s	rtorres	0000005-130116195359883-oozie-oozi-W
17 Jan 2013 01:19:01	KILLED	compla_paraulas	15s	rtorres	0000004-130116195359883-oozie-oozi-W
17 Jan 2013 01:15:27	SUCCEEDED	streaming_wordcount	41s	rtorres	0000003-130116195359883-oozie-oozi-W
17 Jan 2013 01:09:47	SUCCEEDED	streaming_wordcount	45s	rtorres	0000002-130116195359883-oozie-oozi-W
17 Jan 2013 01:07:26	KILLED	streaming_wordcount	20s	rtorres	0000001-130116195359883-oozie-oozi-W
17 Jan 2013 01:05:34	SUCCEEDED	MapReduce	42s	rtorres	0000000-130116195359883-oozie-oozi-W

Figura 6. Quadre de Control d'execució de fluxes.

Els coordinadors (gestionats per ZooKeeper) permeten ordenar l'ordre d'execució de les diferents tasques MapReduce per tal realitzar els tractaments necessaris sobre les dades que han de ser manipulades.

