

- OKKAM
Enabling the Web of Entities

Large-Scale Integrating Project (GA#215032)

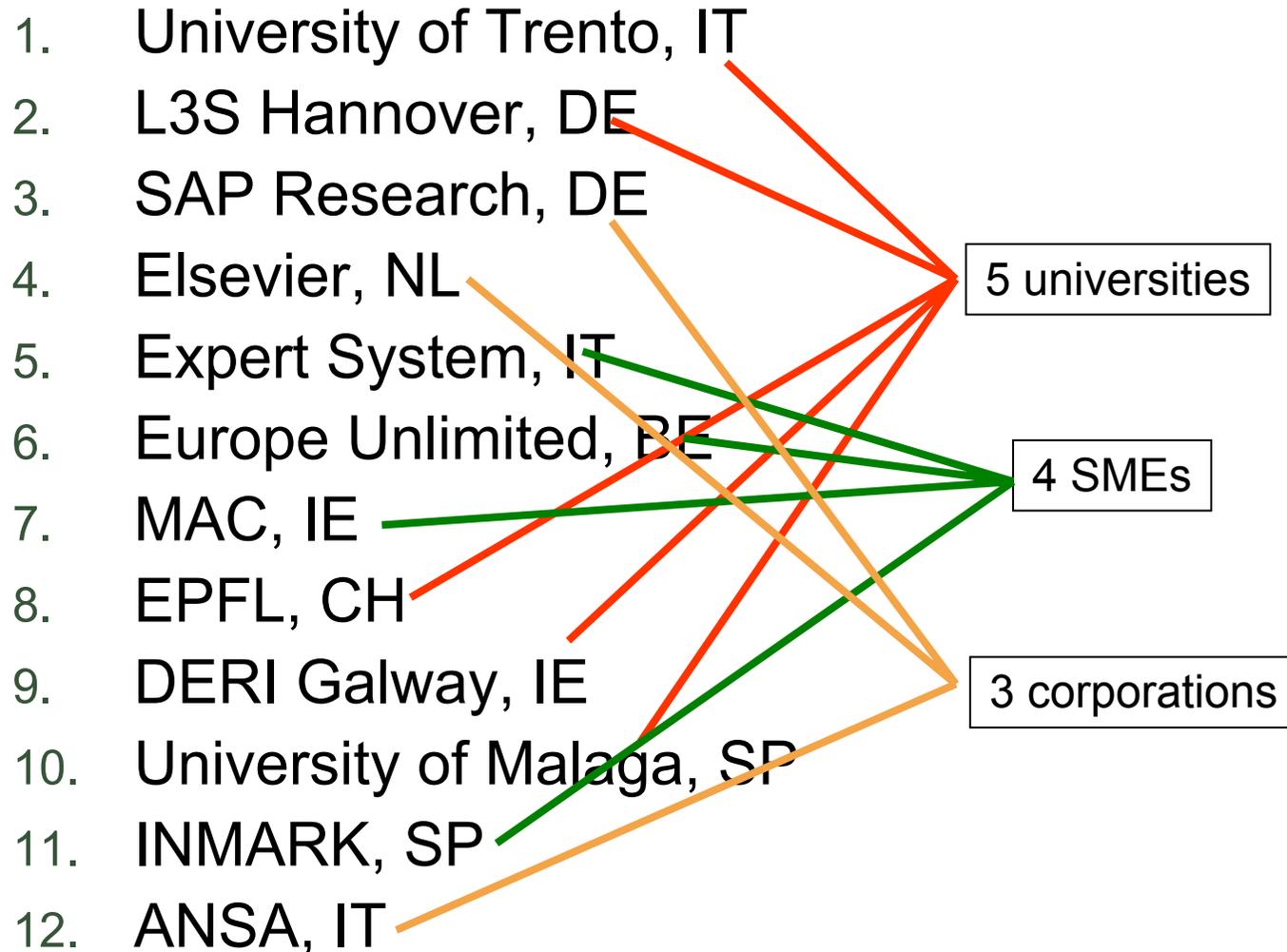


The OKKAM factsheet



<i>Project Title</i>	Enabling the Web of Entities. A scalable and sustainable solution for systematic and global identifier reuse in decentralized information environments
<i>Acronym</i>	OKKAM
<i>Starting date</i>	01/01/2008
<i>End date</i>	30/06/2010
<i>Duration</i>	30 months
<i>No. of partners</i>	12 (5 universities, 4 SMEs, 3 large companies)
<i>Effort</i>	763 person months (more that 63 person years)
<i>Total cost</i>	7.352.931,34 Euro
<i>EU contribution</i>	5.125.000,00 Euro
<i>Type of instrument</i>	Large-scale Integrating Project

The OKKAM consortium



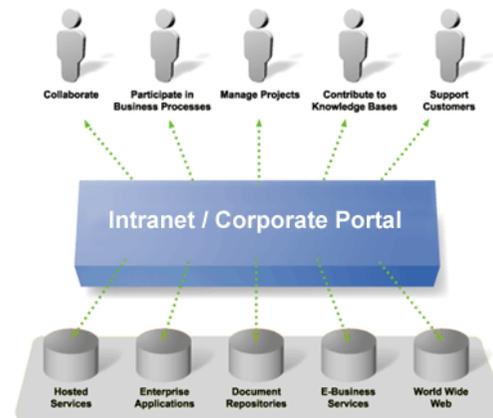
OKKAM: the rational

How many “entities” (persons, locations, organizations, events, projects, products, ...) are named in:

Files in your laptop



Content in your Intranet/
Enterprise Information System

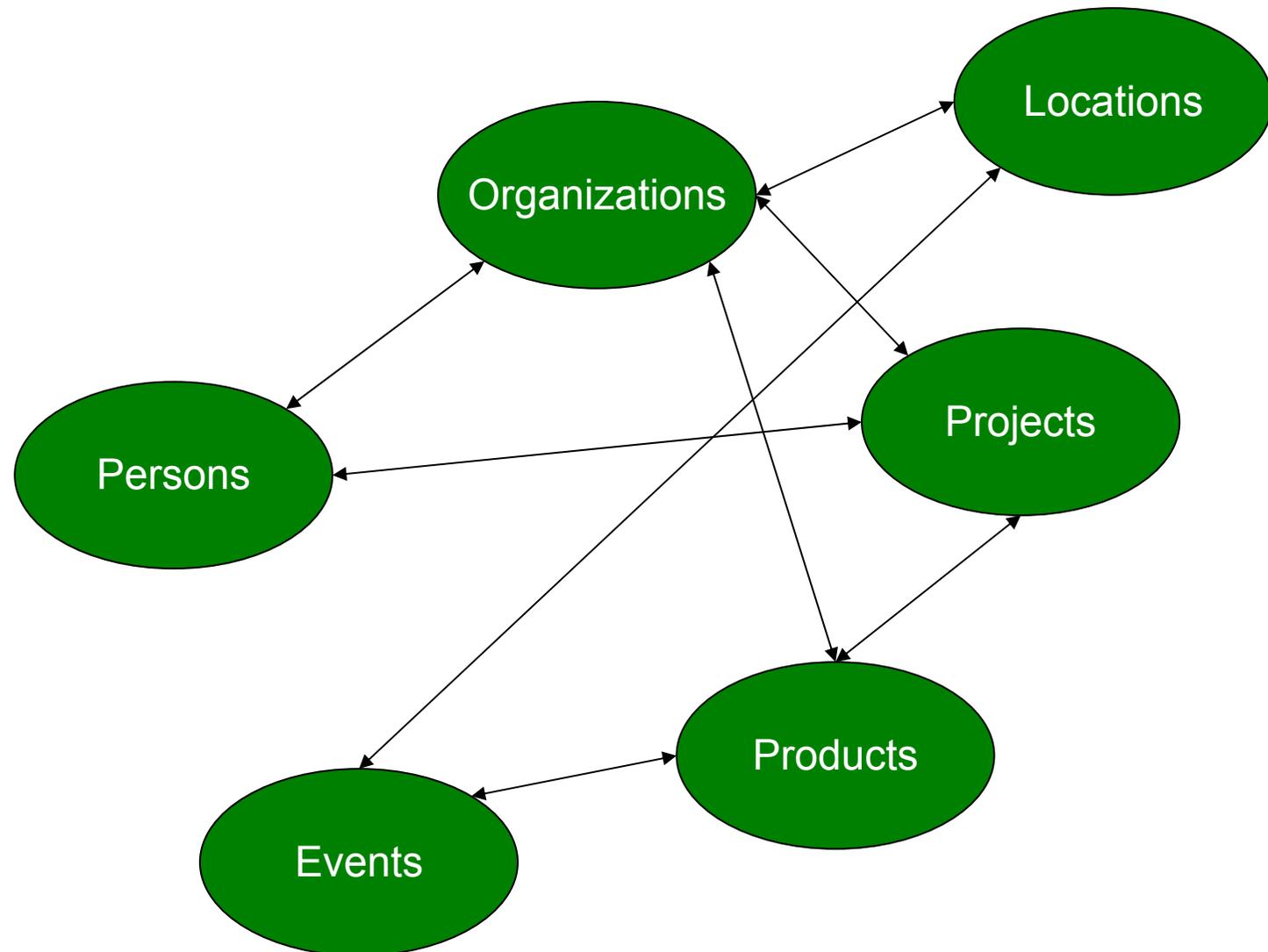


Contents/applications
on the Web



An invaluable asset

“Entities” is what a large part of our knowledge is about:



... or ...

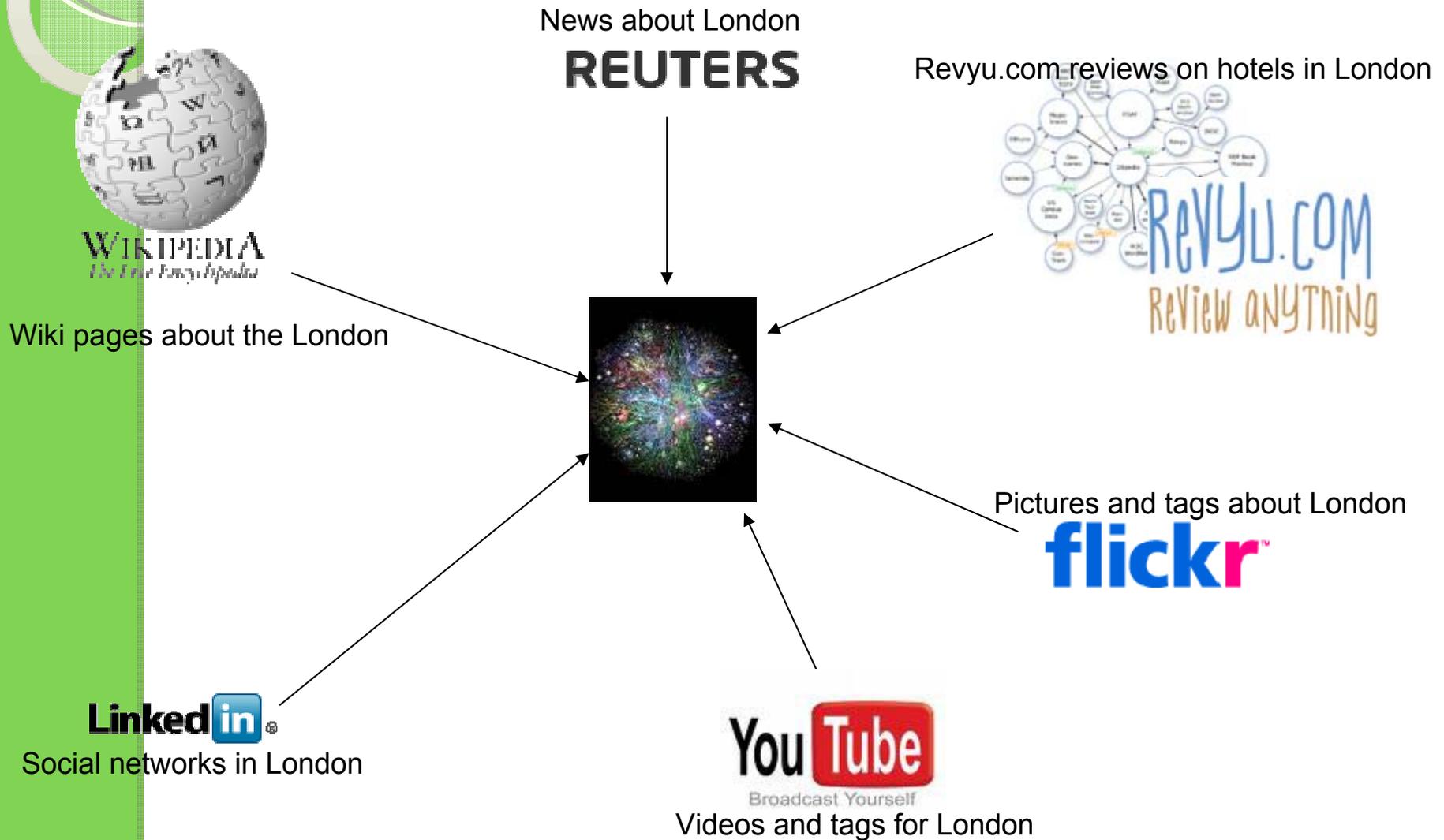
How many “entities” have the same name?

London, KY
London, Laurel, KY
London, OH
London, Madison, OH
London, AR
London, Pope, AR
London, TX
London, Kimble, TX
London, MO
London, MO
London, London, MI
London, London, Monroe, MI
London, Uninc Conecuh County, AL
London, Uninc Conecuh County, Conecuh, AL
London, Uninc Shelby County, IN
London, Uninc Shelby County, Shelby, IN
London, Deerfield, WI
London, Deerfield, Dane, WI
London, Uninc Freeborn County, MN
...

- Or
 - London, Jack
2612 Almes Dr
Montgomery, AL
(334) 272-7005
 - London, Jack R
2511 Winchester Rd
Montgomery, AL 36106-3327
(334) 272-7005
 - London, Jack
1222 Whitetail Trl
Van Buren, AR 72956-7368
(479) 474-4136
 - London, Jack
7400 Vista Del Mar Ave
La Jolla, CA 92037-4954
(858) 456-1850
 - ...

... or ...

How many content types / applications provide valuable information about each of these “entities”?



This is an immense loss of value

- Precision and recall of keyword-based **search engines** is deeply affected
- **Semantic search** still mainly relies on keyword search/matching (as very few URIs are consistently reused across RDF datasets)
- **Information integration** (e.g. from heterogeneous databases) is hard to achieve
- **Mash-ups** from different sources are not easy to create
- **Data/Text/Web mining** tools must struggle with object consolidation from different data sources
- **Information extraction** produces poorly integrated results (co-reference is supported mainly in the same document/corpus, not across large-scale distributed environments)

OKKAM envisions an open, decentralized and global knowledge space in which:

- Every entity (individual, instance, “thing”) is assigned a global identifier, ideally unique
- The same identifier is used to name the same entity across any type of content, format, application, domain, language, culture
- Contents on the Web (from text to multimedia) are consistently annotated with entity unique identifiers
- Users and application are provided with easy and straightforward ways (e.g. GUIs, APIs) for retrieving the IDs of the entities they need to name in their content

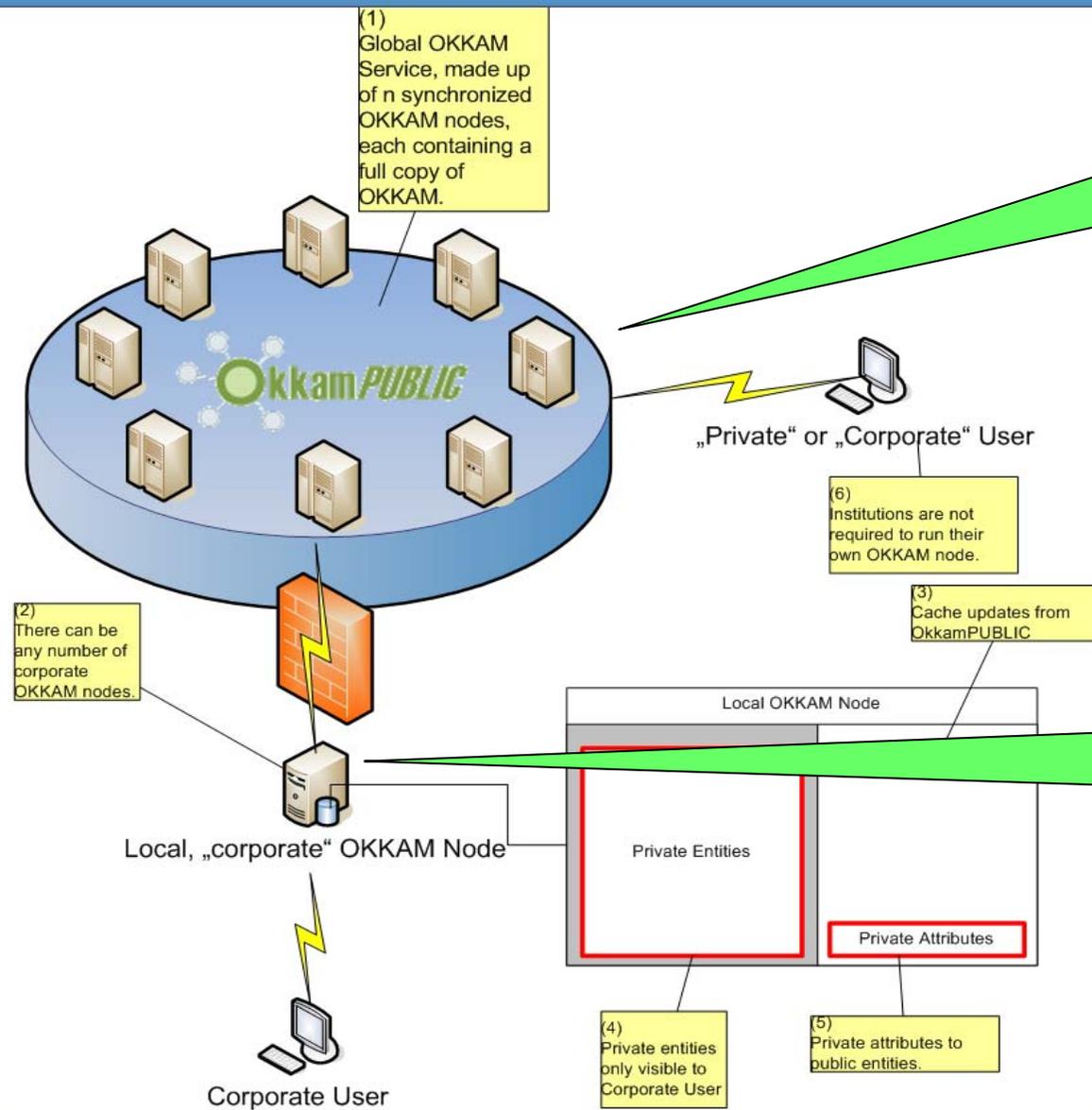
This will enable a new powerful generation of **entity-centric applications and services**

OKKAM aims at implementing an **Entity Name System** (ENS) for the Web, namely an *infrastructure* which can offer the following basic *services*:

- **ID storage and management:** stores, maintains and makes available for reuse IDs (URIs) for anything which is named in a networked environment
- **Entity matching (and look-up):** maps any arbitrary description of an entity to its global ID (or to the collection of known IDs)
- **APIs:** provides application interfaces which can return IDs to applications which needs them in the creation of new content (e.g. word processors, ontology editors, HTML editors, web-based authoring environments, ...)
- **Access:** allow users or application to create new IDs for things which do not have one, or to modify/update an entity description in the ENS

OKKAM ENS – Global and Decentralized

OKKAM Global Distributed Architecture
Working Draft 0.4 of 2008-04-24



Replicated public nodes (ENS servers) hosted in different locations

Local “corporate” nodes for non-public data (and cache)

ENS-empowered tools

Any application which is used to create/edit content can be extended with “plugins” which enable the application to get from the ENS the right ID for a named entity

Examples of OKKAM-empowered tools:



Annotating texts created with MS Word with IDs automatically obtained from the ENS



Reusing URIs in RDF/OWL knowledge bases built with Protégé with IDs automatically obtained from the ENS



Introducing global IDs to content created through web-based authoring interfaces (e.g. FOAF content)

Three OKKAM applications



In the project, we have three application scenarios:

- Entity-centric **semantic search engine** (DERI Galway)
- Entity-centric **organizational knowledge management systems**: entity-centric product lifecycle management within SAP
- Multimedia **authoring environments**: supporting the creation of news (ANSA) and scientific papers (Elsevier) in an ENS-empowered environment

- **Work Packages (Work breakdown structure)**
- **Project Competence Areas (PCAs)**
- **Integration and Demonstration Subprojects (ISP)**
- **Entity-centric Applications**

15 Work Packages

- Further broken down into **Tasks**
- Result in **Deliverables**
 - **Integrated SW Components**
 - **Documents (published, unpublished)**
- Are **not** mapped 1:1 to a single partner
 - but are in the **responsibility** of a single partner
 - require task- or subtask-level **co-operation**

- Bundle related RTD **competences**
- Co-ordinates between **related activities**
- Fosters **coherent** development of contributions
- Ensures **transition** of research results into implementation

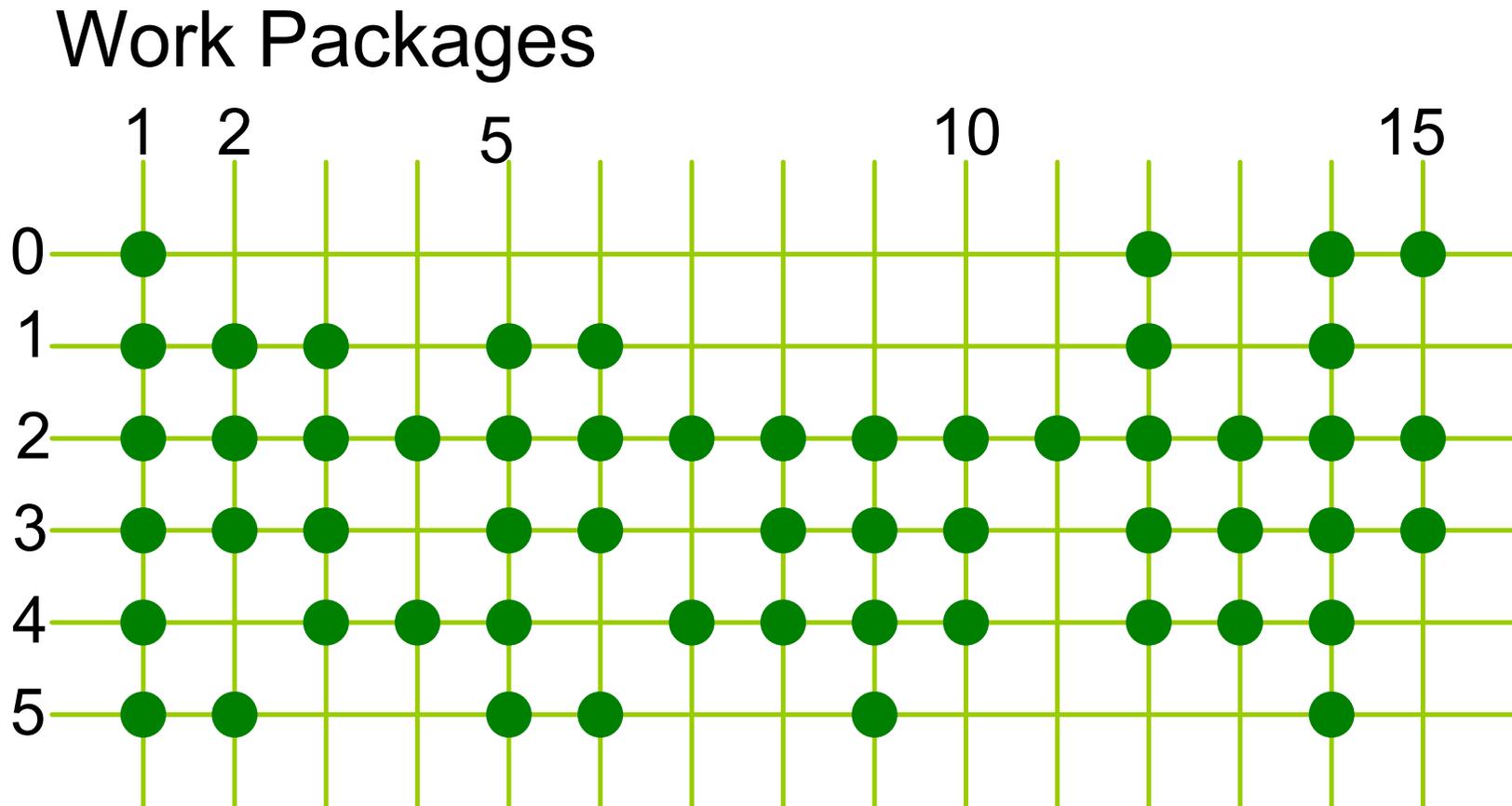
- Are a **crucial instrument for the tight integration** among partners and activities

Project Competence Areas

	Project Competence Area (PCA)	Area Leader
PCA-0	Community building, sustainability and exploitation	Europe Unlimited
PCA-1	Distributed storage & scalable entity management	EPFL
PCA-2	Entity representation, matching and reasoning	Univ. Trento
PCA-3	Entity lifecycle management and evolution	L3S Hannover
PCA-4	Information extraction, crawling and harvesting	ExpertSystem
PCA-5	Security, trust and privacy	Univ. Malaga

PCA vs. Work Packages

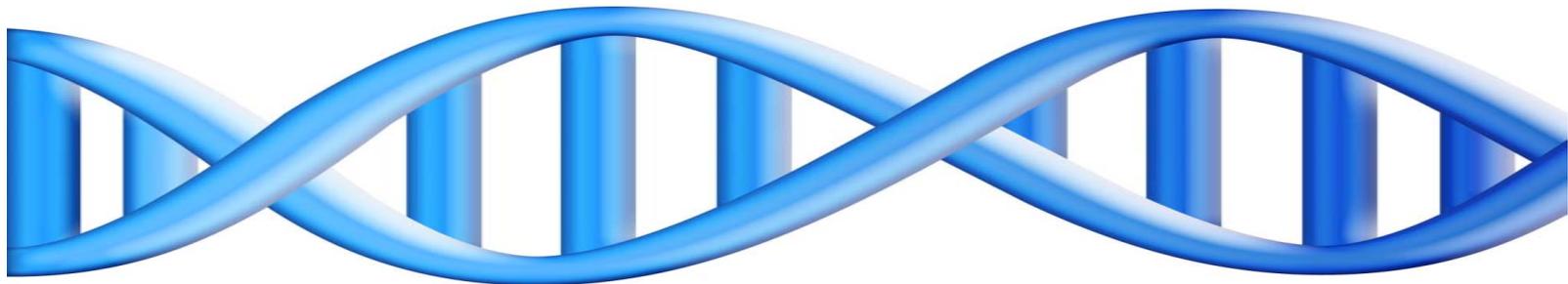
PCAs



- Foster **systematic evolution** of project results
- Ensure **horizontal integration** of RTD contributions
- Result in integration **prototypes**

- ISPs are the **main drivers** during the project life-cycle

- They cause all work to **converge** (at least) 3 times in 30 months



OKKAM Infrastructure Evolution

