# FALCON

# **The Localization Web:**
From L10n workflows to Linked Data

**Dave Lewis**
**CNGL at Trinity College Dublin**

FALCON

- Multilingual web pages could offer an important language resource,
  - e.g. as parallel text for machine translation engine or multilingual term extraction
- Difficult to leverage, HTML is a publication format, it hides valuable translation info:
  - Translated sentence alignment
  - Term meta-data
  - Translation provenance: was it machine translated, transcreated, quality checked?
- Barrier to leverage by industry's long tail of SME LSPs and clients
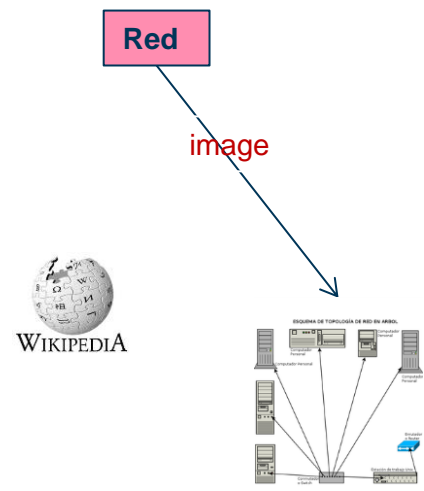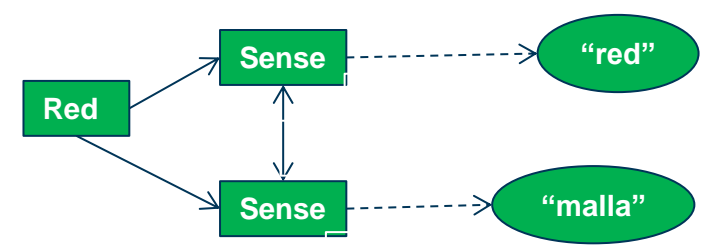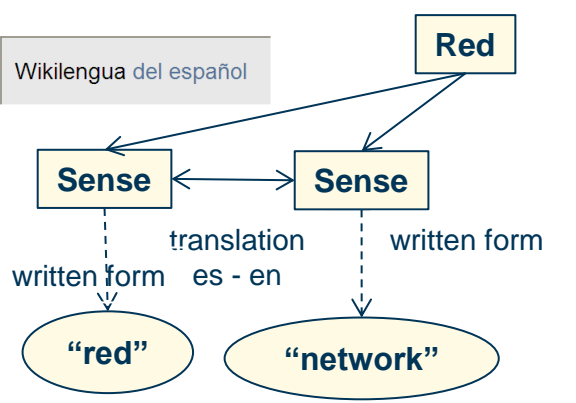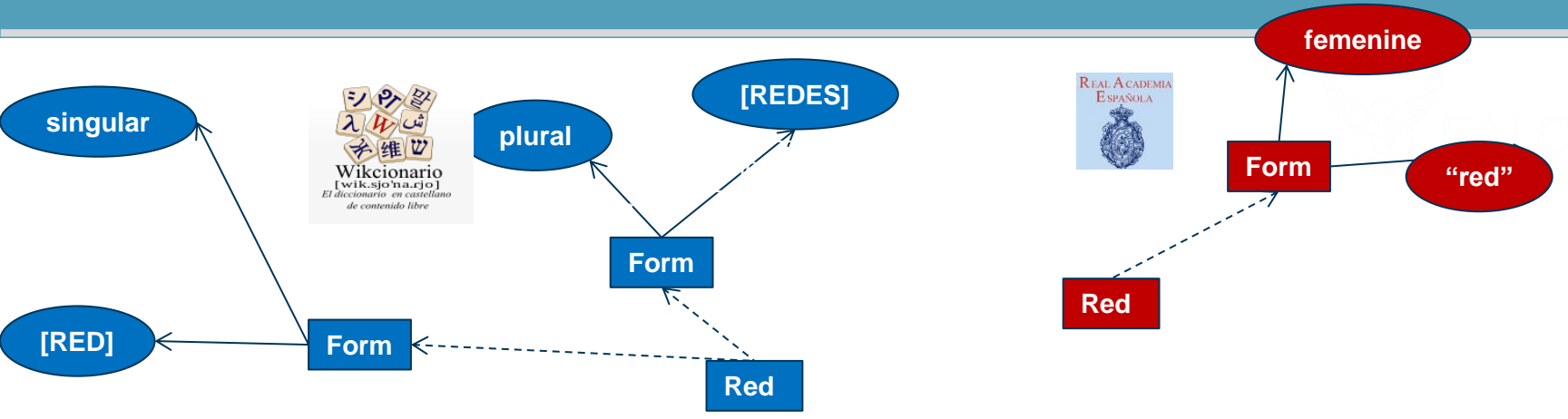
# The Localization Web

- W3C Semantic Web standards allow <u>data</u> to be published on Web
  - Fine-grained URI-based inter-linking
  - Extensible meta-data
  - Standard Query APIs
- Enables a Localization Web
  - Terms and translations become <u>linkable resources</u>
  - Meta-data from L10n workflows <u>adds value</u>
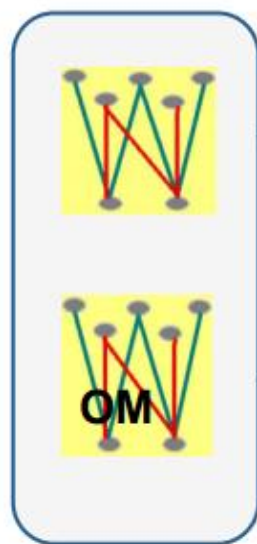  - Leverage in <u>training</u> Machine Translation and Text Analytics

**The Localization Web = Decentralised Annotated
Global Translation Memory and Term Base**

# BabelNet

## Traditional lexical resources

## Collaborative lexical resources



**BabelNet**

A very large multilingual encyclopedic dictionary and semantic network

- fully-structured
- manually curated by experts
- available for a few languages
- difficult to maintain and update

- semi-structured
- collaboratively built by the crowd
- highly multilingual
- up-to-date

# Words as Resources on the Web



FALCON

## The Web of Content

http:// www.ex.org/obama_en.html

Barak Obama if the 44th president of the United State of America. He was first elected in 2009.

Barak Obama si el 44 º presidente de los Estados Unidos de América. Ha fue electo primera vez en 2009.

http:// www.ex.org/obama_es.html

## The Localization Web

### Translation Data

http://data.ex.org/String_0001

**Text:** "Barak Obama if the 44th president of the United State of America."
**Lang:**en

Derived From

Translated From

**Text:**"Barak Obama si el 44 º presidente de los Estados Unidos de América."
**Lang:**es
**TranslatedBy:**Google Translate

Derived From

http:// data.ex.org/String_0002

### Terminology Data

http:// babelnet.org/345621

**Term:** "United State of America."
**Lang:**en

Translation Of

**Term:**"Estados Unidos de América."
**Lang:**es

http:// babelnet.org/57835
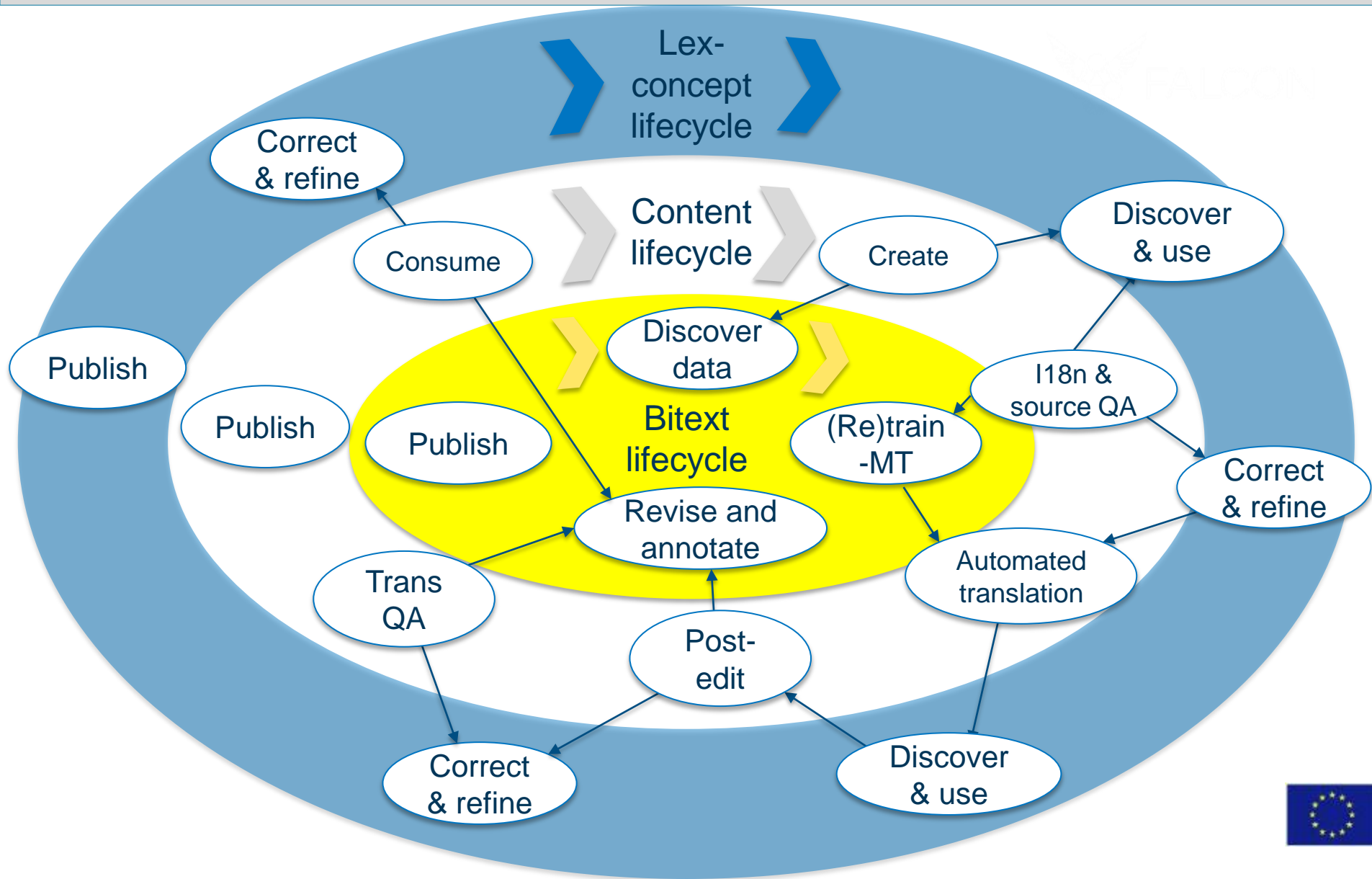
### Encyclopaedic Data

**Topic:** Barak Obama
**Lang:** en
**BirthDate:** 1961-08-04
**Spouse:** Michelle Obama
**Residence:** White House

http:// Dbpedia.org/Page/Barak_Obama

- ## Source Internationalisation
  - – <u>Term disambiguation</u>
  - – <u>Term extraction</u> with translation discovery
  - – <u>Auto-tag named entities</u> with encyclopaedic reference for authors and translators

- ## Machine Translation
  - – <u>Consistent</u> machine translation of terms
  - – <u>Pooling and discovery</u> of parallel text for training

- ## Translation and Post-editing
  - – <u>Term definitions</u> from open encyclopaedic data
  - – <u>Concordancing</u> over a global TM

# Data Management Lifecycles

FALCON

- How can Language Resource Publishers can <u>audit links</u> to and use of resources & <u>track ROI?</u>

- Can Tool Vendors and Integrators <u>expand markets</u> with more open <u>asset management</u> offerings?

- How can SME LSPs gain <u>resource sharing and pooling</u> opportunities and avoid lock-in?

- Can LSPs and clients can use Active Curation to quickly <u>train domain specific</u> SMT and text analytics components?

# Data Management Needs?

- Assert ownership and attribution, licensing, access control?

- Persistent URLs?

- Open royalty free standards?

- Indexing is key, federated vs aggregated data?

- Third party submission of errors, QA, corrections? Publish action on submissions?

- Query bitext on: languages, terms, MT engine, MT confidence, QA, translator qualfications, postedit characteristics?

- Query term base on: language, domains, context, semantic relations, provenance of lexical/conceptual data?

- Web API:
  - HTTP content negotiation (unicode extensions for translation?),
  - Format: RDF, TMX, TBX, RDF, JSON-LD,
  - SPARQL?

# More Information

- Contact: dave.lewis@cs.tcd.ie

- http://www.falcon-project.eu


- See also:
  - Linked Data for Language Technology (LD4LT) W3C Community Group
  - http://www.w3.org/community/ld4lt/