

Technological tools for dictionary and corpora building for minority languages: example of the French-based Creoles

Paola Carrión Gonzalez(1,2), Emmanuel Cartier(1)

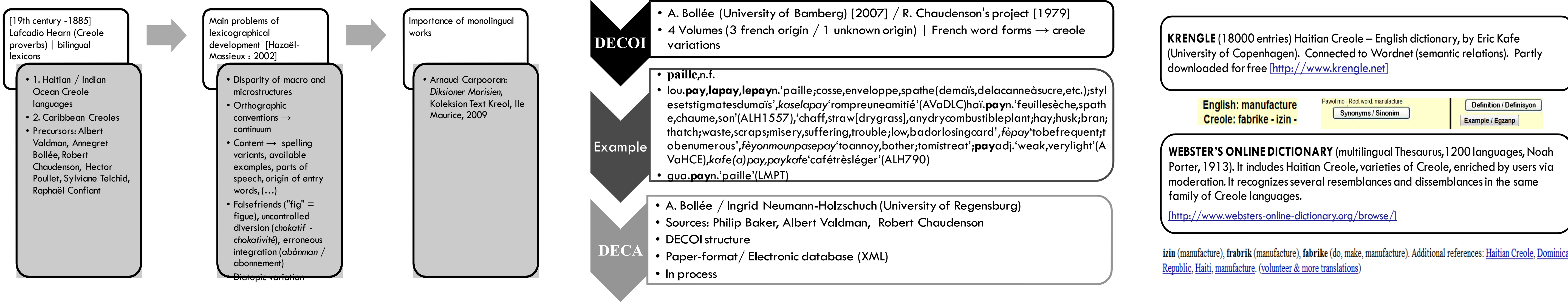
(1)LDI, CNRS UMR 7187, Université Paris 13 PRES Paris Sorbonne Paris Cité

(2)Departamento de Traducción e Interpretación, Facultad de Filosofía Y Letras, Universidad de Alicante

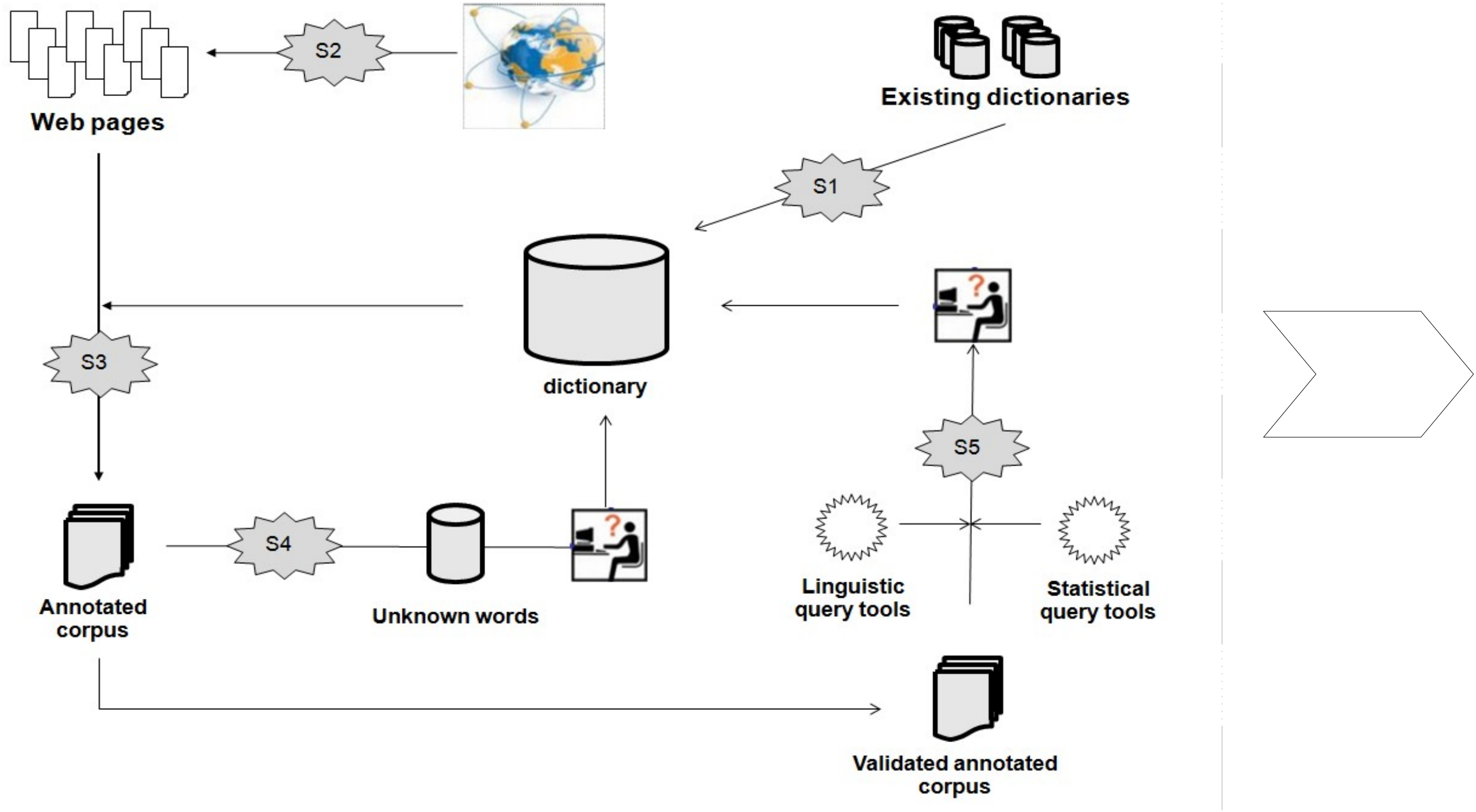
E-mail : pccg1@alu.ua.es , ecartier@ldi.univ-paris13.fr

We present a project which aims at building and maintaining a lexicographical resource of contemporary French-based creoles, still considered as minority languages, especially those situated in American-Caribbean zones. These objectives are achieved through three main steps: 1) Compilation of existing lexicographical resources (lexicons and digitized dictionaries, available on the Internet); 2) Constitution of a corpus in Creole languages with literary, educational and journalistic documents, some of them retrieved automatically with web spiders; 3) Dictionary maintenance: through automatic morphosyntactic analysis of the corpus and determination of the frequency of unknown words. Those unknown words will help us to improve the database by searching relevant lexical resources that we had not included before. This final task could be done iteratively in order to complete the database and show language variations within the same Creole-speaking community. Practical results of this work will consist in 1/ A lexicographical database, explicating variations in French-based creoles, as well as helping normalizing the written form of this language; 2/ An annotated corpora that could be used for further linguistic research and NLP applications.

1. Existing Lexicographical Resources in French-based Creole



Project Architecture



This architecture comprises five main steps:

- Step 1 (S1): this preliminary step aims at building a first electronic compilation from existing lexicographical resources. This implies gathering the existing resources, either digitized resources from paper-based dictionaries, or fully electronic resources. This also implies identifying available resources. This initial step is detailed in section 4;
- Step 2 (S2) : this second step aims at setting up an infrastructure enabling corpora building and feeding; this means identifying web available resources as well as setting up automatic procedures to retrieve on a regular basis these documents; it is detailed in section 5.1.
- Steps 3 and 4: (S3-S4) morphosyntactic analysis of the corpora to maintain the existing dictionary; automatic analysis of corpora will explicit unknown words, and some of them will have to be included in the initial dictionary; iteration of this procedure will permit to complete the existing dictionary, as well as enabling morphosyntactic annotation of the corpora; this step is detailed in section 5.2.
- Step 5 (S5) : this step, out of the scope of this paper, will be implemented as soon as the dictionary is sufficiently completed; annotated corpus could then be validated and then queried using linguistic and statistical tools, so as to improve information in the dictionary. It will be evoked in the section 5.3.

2. Dictionary building: existing resources compilation

Objectives :

- retrieve available (free) data from the web
- two main sources : paper-based and electronically-based dict.

Steps followed :

- compilation of available resources
- scripts to retrieve data and combine them into a unique dictionary keeping track of source and

Electronically-based dict.: about 2000 entries

- purpose of the dictionaries : mainly educational or popularisation
- small lexicons
- weaknesses of linguistic information (mostly reduced to entry and definition / translation)

Paper-based dict.: about 20 000 entries

- less numerous
- contain much more linguistic information (pos, examples, translation, spelling variants, cross-references, semantic relations, etymology)
- generally bilingual resources
- require more complex data processing.

Web Corpora can help building/maintaining lexicographical resources

Our project will build its dictionary not only from existing electronic resources, but also by morphosyntactically analyzing corpora (with the help of the previously setup dictionary) (Kilgariff, 2003; Baroni, 2009 for example; and for minority languages Scannel, 2007)

Steps :

- corpora downloading (on a regular basis),
- automatic morphosyntactic analysis,
- manual dictionary improvement.

Iterative process :

- 1/ extracting unknown words as a first step
 - 2/ checking meanings associated to recognized words
- => will end up with a more exhaustive dictionary, following the Zipf's law.

Three main corpora:

- first : 1 million words, setup to tune the morphosyntactical analysis and the dictionary improvement process, with in mind the general-purpose language;
- second : the Bible, to complete and tune the dictionary with a specific vocabulary;
- third : evolving corpus, to maintain the dictionary and setup an adequate environment for lexicographers.

Analysis :

1/Quick lexicographic coverage of the dictionary: from 28,6% unknown words with about 70% unique words (first corpus) to 7,6% and about 98% of unique words (second corpus)

=> whereas a relatively small corpus enable to cover 90% of lexicographic entries, only a really huge corpus enable to tend to 100% coverage. (Zipf's law)

2/ Dictionary coverage : 123 245 unique words-forms; coverage to be tuned with :

- **phrases** (about 50% of the vocabulary, see Sag et al, 2002, for example);
- **unknown words without linguistic information**
- **known words linguistic information**
- **spelling variants.**

3/ unknown words

it would be necessary to have automatic procedures to track the meaning evolutions, rather than only word forms existence.

- words from other language (specifically from English, in our case),
- misspelled words,
- proper names,
- specific notations,
- real unknown words.

=> The resulting list has been included in the dictionary without any linguistic information, except for a small part of it with information taken from web-based dictionary (that could not be retrieved globally, but can be used for individual word search) or paper-based dictionaries (see Carrión, 2011 for the list of these dictionaries). For each of these words, we have decided, in a first step, to include only part of speech, translation in French and English, and varieties if applicable.

4. Conclusion

On-going project whose goals are to:

- 1/ **Explicit a methodology** to improve NLP and linguistics research for minority-languages, focusing on the American-Caribbean French-Based Creoles;
- 2/ **Setup procedures and an environment for lexicographers** to store and maintain lexicographical data from existing resources and web corpora.

Results / Perspectives

- 1/ **Compilation of existing lexicographical resources** : about 20 000 lexicographical entries, but exhibited complex-to-solve problems : quality, quantity diverge from one resource to another, macro and micro-structures are far from unified. => Discussion to work with the DECA team;
- 2/ **Dictionary completion and maintenance through web corpora analysis** :
=> In the course of setting up a web-based environment to maintain the dictionary and study lexicographical phenomena through several iterative processes:
 - Live corpus download (BootCat, RSS Corpus Builder);
 - Cleaning and normalization of web pages;
 - POS tagging with existing lexicographical resources;
 - Web-based environment for lexicographers to study corpus-based lexicological phenomena; (IMS Open Corpus Workbench);

This project has generated two crucial elements for the American-Caribbean creoles: a POS annotated corpus and a POS tagger. These data and tools will be soon released as open-source.

