

**Premier inventaire (non exhaustif) des corpus avec annotation de haut niveau pour le français écrit contemporain disponibles ou interrogeables en ligne**

Assemblée générale de l'Infrastructure de recherche Corpus, Consortium Corpus Ecrits, 23 novembre 2012

Amalia Todirascu (LILPA, Université de Strasbourg), Agnès Tutin (LIDILEM, Université de Grenoble)

N'hésitez à signaler d'autres ressources (écrire à [todiras@unistra.fr](mailto:todiras@unistra.fr) ou [agnes.tutin@u-grenoble3.fr](mailto:agnes.tutin@u-grenoble3.fr))

Site du groupe : <https://listes.cru.fr/wiki/corpus-ecrits/public/groupe-8>

Type d'annotation		Nom	Brève description	Références
Annotations syntaxiques	Corpus annotés	<b>French Treebank</b>	<ul style="list-style-type: none"> <li>• Extraits du Monde</li> <li>• correction manuelle</li> <li>• constituants et fonctions (sujets, objets etc.), fonctions de surface</li> <li>• compatibilité avec plusieurs analyseurs</li> </ul> Disponible pour la recherche sur demande <a href="http://www.llf.cnrs.fr/Gens/Abeille/French-Treebank-fr.php">http://www.llf.cnrs.fr/Gens/Abeille/French-Treebank-fr.php</a>	Abeillé, A., L. Clément, and F. Toussenet. 2003. 'Building a treebank for French', in A. Abeillé (ed) <i>Treebanks</i> , Kluwer, Dordrecht.
		<b>Séquoia</b>	Multi-domaine, multi-genre, <b>libre de droit</b> constituants et dépendances <a href="https://www.rocq.inria.fr/alpage-wiki/tiki-index.php?page=CorpusSequoia">https://www.rocq.inria.fr/alpage-wiki/tiki-index.php?page=CorpusSequoia</a>	Candito M.-H. and Djamé Seddah, 2012, Le corpus Sequoia : annotation syntaxique et exploitation pour l'adaptation d'analyseur par pont lexical, <i>Proceedings of TALN'2012</i> , Grenoble, France.
		<b>Ecritures</b>	corpus constitué de plusieurs états (brouillons) de rapports sociaux rédigés dans le cadre de la protection de l'enfance (300000 mots). Annotation syntaxique ad hoc avec Le Trameur. Projet en cours Disponible sous conditions (corpus sensible) <a href="http://www.univ-paris3.fr/anr-ecritures/">http://www.univ-paris3.fr/anr-ecritures/</a>	<b>CISLARU Georgeta, SITRI Frédérique</b> , "De l'émergence à l'impact social des discours : hétérogénéités d'un corpus", <b>Langages, N°187 (3/2012)</b>
		<b>CLASSYN</b>	Corpus constitué des articles scientifiques et articles de vulgarisation dans les domaines de la médecine et de l'informatique. Annotation syntaxique automatique (analyseur en dépendances B.Bohnet). Disponible sur demande à <a href="mailto:todiras@unistra.fr">todiras@unistra.fr</a>	<b>Todirascu, A., Pado, S., Krisch, J., Kisselew, M. Heid</b> , U. French and German Corpora for Audience-based Text Type Classification, Proceedings of the Eight International Conference on Language Resources and Evaluation (LREC'12)
	Corpus interrogeables en ligne	<b>L'arboratoire</b> (interrogeable en ligne)	Syntaxe de contraintes. Textes français interrogeables en ligne (Europarl) <a href="http://corp.hum.sdu.dk/arboratoire.html">http://corp.hum.sdu.dk/arboratoire.html</a>	Salmon-Alt S., Bick E., Romary L., Pierrel J.M., « La FReeBank : vers une base libre de corpus annotés », <i>Actes de TALN 2004</i> , 18-23 avril 2004.
		<b>Scientext</b> (interrogeable en ligne)	(écrits scientifiques interrogeables en ligne) (analyse avec Syntex, syntaxe de dépendance, Bourigault) <a href="http://scientext.msh-alpes.fr/">http://scientext.msh-alpes.fr/</a>	Falaise A., Tutin A., Kraif O. (2012), Une interface pour l'exploitation de corpus arborés par des non-informaticiens : la plate-forme ScienQuest du projet Scientext, <i>Traitement Automatique des Langues</i> 2011

					Volume 52 Numéro 3, <a href="http://www.atala.org/IMG/pdf/Falaise-TAL52-3.pdf">http://www.atala.org/IMG/pdf/Falaise-TAL52-3.pdf</a>
Annotations sémantiques	Corpus annotés	Entités nommées	<b>Etape Ester2</b>	Corpus utilisé dans des campagnes d'évaluation des campagnes d'évaluation, annotation des personnes, lieux, organisation (transcription d'émissions de radios)  <b>Diffusé par ELRA</b>	Sylvain Galliano, Edouard Geoffrois, Djamel Mostefa, Khalid Choukri, Jean-François Bonastre, Guillaume Gravier: The ESTER phase II evaluation campaign for the rich transcription of French broadcast news. <i>INTERSPEECH 2005</i> : 1149-1152 <a href="http://www.afcp-parole.org/etape/docs/interspeech05-ester.pdf">http://www.afcp-parole.org/etape/docs/interspeech05-ester.pdf</a>
		Sentiments et opinions	<b>Blogoscopie</b>	Blogoscopie (ANR) (piloté par le LINA) : annotation des blogs avec les concepts et les évaluations. <b>Téléchargeable après accord</b> <a href="http://www.lina.univ-nantes.fr/?Blogoscopie,762.html">http://www.lina.univ-nantes.fr/?Blogoscopie,762.html</a>	DAILLE B., DUBREIL E., MONCEAUX L. ET VERNIER M. Annotating opinion—evaluation of blogs : the Blogoscopy corpus. <i>Language Resources and Evaluation</i> Springer 2011.
		Expressions temporelles	<b>French TimeBank</b>	Repérage des entités temporelles et des relations entre entités dans un corpus de textes journalistiques (Est Republicain) <b>Librement téléchargeable</b> <a href="https://gforge.inria.fr/projects/fr-timebank/">https://gforge.inria.fr/projects/fr-timebank/</a>	Bittar A., Amsili P., Denis P. (2011), French TimeBank : un corpus de référence sur la temporalité en français, <i>TALN 2011 - Traitement Automatique des Langues Naturelles</i> , (2011) 259-270 <a href="http://www.linguist.univ-paris-diderot.fr/~amsili/papers/TALN-2011-French-TimeBank.pdf">http://www.linguist.univ-paris-diderot.fr/~amsili/papers/TALN-2011-French-TimeBank.pdf</a>
		Désambiguïsation	<b>Romanseval</b>	Corpus de la campagne ROMANSEVAL (60 mots désambiguïsés).  <b>Diffusé par ELRA.</b>	<a href="http://aune.lpl.univ-aix.fr/projects/romanseval/">http://aune.lpl.univ-aix.fr/projects/romanseval/</a>
	Corpus interrogeable en ligne	Nominalisation	<b>Nomage</b>	Corpus avec annotation des nominalisations (French Treebank)  <a href="http://nomage.recherche.univ-lille3.fr/nomage/">http://nomage.recherche.univ-lille3.fr/nomage/</a>	Balvet A., Barque L., Condet M.-H., Haas P., Huyghe R., Marín R., Merlo A. (2012). La ressource Nomage, Confronter les attentes théoriques aux observations du comportement linguistique des nominalisations en corpus, <i>TAL</i> , Volume 52 – n° 3/2012, pp. 1-24. <a href="http://www.atala.org/IMG/pdf/Balvet-TAL52-3.pdf">http://www.atala.org/IMG/pdf/Balvet-TAL52-3.pdf</a>
Annotations textuelles et discursives	Corpus annotés	Annotations textuelles	<b>Annodis</b>	Relations de discours entre les unités minimales de discours, macro-structures, notamment les structures énumératives, chaînes topicalisées. Corpus variés. <b>Librement disponible</b> : <a href="http://redac.univ-tlse2.fr/corpus/annodis/">http://redac.univ-tlse2.fr/corpus/annodis/</a>	Péry-Woodley M.-P., Afantenos S. D., Ho-Dac L.-M., Asher N. (2011). La ressource ANNODIS, un corpus enrichi d'annotations discursives. <i>TAL 52(3)</i> , pp 71-101. <a href="http://redac.univ-tlse2.fr/corpus/annodis/biblio/peryEtAl2011-TA52-3.pdf">http://redac.univ-tlse2.fr/corpus/annodis/biblio/peryEtAl2011-TA52-3.pdf</a>

		<b>Géopo</b>	Corpus de 270 000 mots d'articles expositifs autour des relations internationales. Structure textuelles, sections, énumérations. <b>Librement disponible :</b> <a href="http://redac.univ-tlse2.fr/corpus/geopo.html">http://redac.univ-tlse2.fr/corpus/geopo.html</a>	Lydia-Mai Ho-Dac : <i>La position initiale dans l'organisation du discours : une exploration en corpus</i> , thèse de doctorat, Université de Toulouse-Le Mirail, novembre 2007.
Annotations discursives		<b>Lelie</b>	Corpus de textes du genre procédural (procédures, exigences, textes didactiques) 8000 textes procéduraux en français et environ 1500 textes procéduraux en anglais, 2000 exigences français et 1200 en anglais. Annotations des relations du discours, structure des verbes  Adresse du projet : <a href="http://www.irit.fr/recherches/ILPL/lelie/accueil.html">http://www.irit.fr/recherches/ILPL/lelie/accueil.html</a>	
		<b>French Discourse Treebank</b>	Objectif : construire un corpus similaire à Penn Discourse Treebank. Analyse des connecteurs inspirée de SDRT et RST (en cours). Projet porté par ALPAGE	
		<b>Projet COMTIS</b>	Extraits d'Europarl annotés avec connecteurs de discours <a href="http://www.idiap.ch/dataset/europarl-direct/">www.idiap.ch/dataset/europarl-direct/</a>	Popescu-Belis A., Meyer T., Liyanapathirana J., Cartoni B. & Zufferey S. 2012. <i>Discourse-level Annotation over Europarl for Machine Translation: Connectives and Pronouns</i> . Proceedings of LREC 2012, May 23-25 2012, Istanbul, Turkey.
Coréférence et anaphores		<b>Corpus d'expressions anaphoriques grammaticales</b>	Expressions anaphoriques grammaticales (presse, écrits scientifiques, ouvrages scientifiques). 1 million de mots  <b>Diffusé par ELRA</b>	Tutin A., Trouilleux F., Clouzot C., Gaussier E., Zaenen A., Rayot S., Antoniadis G. (2000). Annotating a large corpus with anaphoric links. <i>Proceedings of DAARC 2000</i> . Lancaster 16-18 November 2000 <a href="http://halshs.archives-ouvertes.fr/docs/00/37/33/27/PDF/tutin_daarc2000.pdf">http://halshs.archives-ouvertes.fr/docs/00/37/33/27/PDF/tutin_daarc2000.pdf</a>
		<b>Dédé</b>	Corpus de descriptions définies. <b>Librement disponible :</b> <a href="http://www.cnrtl.fr/corpus/dede/">http://www.cnrtl.fr/corpus/dede/</a>	Gardent C., Manuélian H. (2005). Création d'un corpus annoté pour le traitement des descriptions définies. <i>TAL</i> . 2, pp 1-15 <a href="http://www.loria.fr/~gardent/publis/tal05-dede.pdf">http://www.loria.fr/~gardent/publis/tal05-dede.pdf</a>
		<b>EvalRefGen</b>	Corpus annoté en chaînes de référence (multigenre) <b>Librement disponible</b> sur demande à l'auteur <a href="mailto:todiras@unistra.fr">todiras@unistra.fr</a>	Laurence Longo, Amalia Todirascu: RefGen: a Tool for Reference Chains Identification. IMCSIT 2010: 447-454
		<b>Projet MC4</b>	Etude de chaînes de coréférence dans des textes du français médiéval et contemporain, dans les textes narratifs et non-narratifs. Projet en cours.	