

英語コーパス学会第48回大会

2022年10月1日 (土) 9時30分~17時40分 On Zoom

Workshop: Data-Driven Course Design

Tom Cobb

Université du Québec à Montréal

cobb.tom.3@gmail.com

Sat Oct 1 2022 9:40-10:40



Running a Vocabulary Course With Lextutor

Dr. Tom Cobb (Université du Québec à Montréal)

The importance of vocabulary knowledge in any type of language course or curriculum is now acknowledged but still not easy to incorporate in a systematic manner. There are few dedicated vocabulary courses, and the ones there are do not match the increasing specialisation of many learner programs. Lextutor has been designed to basically take on the whole job of a dedicated vocabulary supplement from a practical corpus perspective. My workshop will show the main steps in this process, from building a corpus of learning materials, to placement testing, to text selection and adaptation, to assuring a manageable supply of new and appropriate items, to test writing that reflects what learners have actually been exposed to sufficiently to be tested on. The theme is 'corpora for courses' and I will share results from locales where this approach and technology is being deployed.

Premise and RQ

Much learning research has been done

And replicated

Many good corpus tools based on it are available

And online or free to download
 How often/systematically are these be brought together
 for the benefit of learners?

• IE, as courses? Curricula?

What would doing so look like? SOME CONCRETE IDEAS ->

Definitions, context, rationale

Corpus

- Representative collection of texts
 - Ex Brown Corpus (1960s)
 - 1 million words
 - 500 textes on 15 topics of 2000 words each
 - BNC (1980s)
 - 100 million words
 - >4000 texts on 100 topics of 25,000 words (approx.)
 - COCA, Subtlex, TenTen ...

Corpus – read and interpreted how?

```
001. 0 décembre, 134 personnes 53 hommes, 55 FEMMES 26 enfants de un à treize ans ét
002. 

de ces deux affaires, une homme et une FEMME : Alfred Sirven, ancien directeur
003. 
ns cent quatre-vingt-dix hommes et deux FEMMES ? Un peu plus tôt, a-t-il partici
004. Description chorégraphie. La pièce met en scène une FEMME agressée par deux hommes et appela
005. pide peur de l'homme, des hommes et des FEMMES aussi... Un besoin maladif d'igno
006. août 1942. Déjà, de nombreux hommes et FEMMES avaient été astreints à des pério
007. amenant vers l'emploi des hommes et des FEMMES dans la force de l'âge, souvent c
008. Un'hématocrite chez 46 athlètes hommes et FEMMES de haut niveau (participant à des
009. 

Les hommes d'aujourd'hui traitent les FEMMES de manière aussi primitive, aussi
010. 

de vingt-cing ans et aux hommes et aux FEMMES de plus de vingt-cing ans. En rev
011. Iong desguelles se succèdent hommes et FEMMES de toutes conditions et origines
012. 

ne se reproduisent, tous les <u>hommes</u> et <u>FEMMES</u> du groupe qui travaillent à la re
013. Description hommes d'affaires tout-puissants... "La FEMME du président" a encore six heures
014. Tavaillait comme un homme, ou comme une FEMME Elle croit ce qu'elle dit et dit
015. Onnage principal, l'homme trompé par sa FEMME employant maints stratagèmes pour
016. ports ambigus d'un vieil homme et d'une FEMME en fleur. Une tragédie intimiste s
017. 

janvier, traite ainsi des hommes et des FEMMES en politique dans le monde, du "m
018. Int ce partage de la vie entre hommes et FEMMES en termes politiques et culturels
019. 
emportent maintenant des hommes et des FEMMES encore jeunes qui laissent des en
020. 

e malade et doit subir la cruauté de sa FEMME Enlevé par ses hommes, il finit p
021. In rieusement atteinte ? "J'espère que les FEMMES et aussi les hommes de ce pays e
022. vant ces ouvrages historiques signés de FEMMES et d'hommes publics, dont on imag
023. U du Soudan ROUX OLIVIER DES hommes, des FEMMES et des enfants meurent de faim au
024. 

le monde a toujours fermé les yeux. Des <u>FEMMES</u> et des <u>hommes</u> n'ont cessé de tire
025. es", "filières criminelles") contre des FEMMES et des hommes prêts à se laisser 4
```

Definitions...

"Data driven learning" or "corpus based learning"

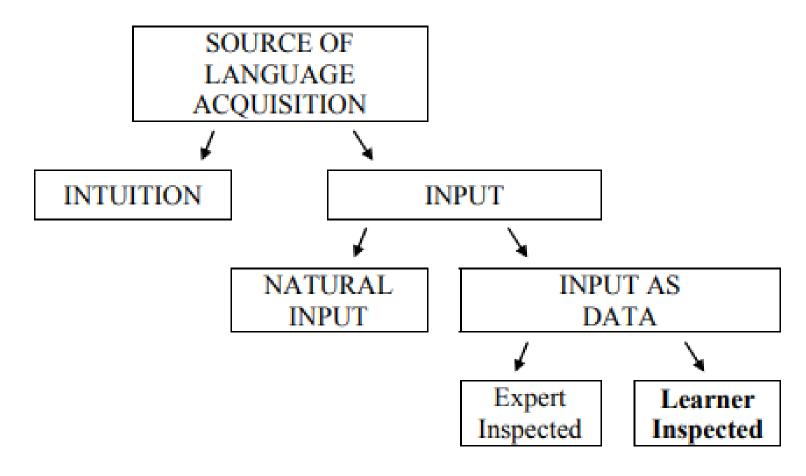


Figure 1. The place of data-driven language learning in the broader scheme

Main idea of DDL

Language learning happens through input,

BUT

- Language input reveals its patterns slowly
 - Too slowly to be effective in most L2 learning
- Input, however, can be treated by computer programs to expose patterns more quickly
 - more clearly
- To create more transparent, 'learnable' input
 - Patterns exposed through data 'agglomeration'
 - As compared to rules or explanation

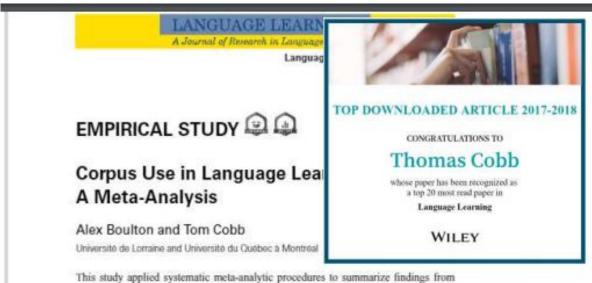
DDL (Data driven learning) is well established as a set of learning tools

- Concordances, frequency lists, learner-as-linguist, discovery learning, pattern extraction, exemplar-based not rule-based learning, etc.
 - But does it work?
- Boulton & Cobb (2017) meta-analysis of DDL
 - Effect size of 1.5 cf 'traditional' learning method
 - Meaning scores 1.5 standard deviations higher
- Many ESL courses now employ these tools piecemeal
 - But can they be integrated systematically?
 - At the point of course conception?

Learning from corpora (=DDL)

Apprendre à travers l'utilisation des corpus / Apprentissage faisant appel aux corpus

• Récemment montré efficace de façon concluante (2017)



This study applied systematic meta-analytic procedures to summarize findings from experimental and quasi-experimental investigations into the effectiveness of using the tools and techniques of corpus linguistics for second language learning or use, here referred to as data-driven learning (DDL). Analysis of 64 separate studies representing 88 unique samples reporting sufficient data indicated that DDL approaches result in large overall effects for both control/experimental group comparisons (d=0.95) and for pre/posttest designs (d=1.50). Further investigation of moderator variables revealed that small effect sizes were generally tied to small sample sizes. Research has barely begun in some key areas, and durability/transfer of learning through delayed posttesting remains an area in need of further investigation. Although DDL research demonstrably improved over the period investigated, further changes in practice and reporting are recommended.

Comparé à une gamme d'autres moyens d'apprendre...

Le vocabulaire

Les locutions

Le syntaxe

La grammaire

La culture

... DDL 'GAGNE' par la taille de l'effet de 1,5 écarts-types (**effect size of 1.5**) pour les études pré-post

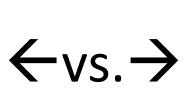
Ex. Control Group Mean 65, SD 15
DDL Group Mean 80, SD 15

Example of a DDL learning activity?

Constructing word meanings from raw data vs pre-constructed

knowledge







Qu'est-ce qu'une activité d'apprentissage basée sur un corpus?

Exemple 2 – sensibilisation à la culture

<u>Tâche</u>

- Do concordances for 'woman' with collocation 'man' on <u>either</u> side in English, French, and Spanish corpora
 - Woman man
 - Femme homme
 - Mujer hombre

Then in a one-page report:

- Is the sequence always man first then woman second across these three languages?
- In which language is there most variation in the sequence?
- Is the common sequence formulaic or indicate a belief in the superiority of one gender?
 - Or with one gender acting and the other being acted upon?

```
001. 0 décembre, 134 personnes 53 hommes, 55 FEMMES 26 enfants de un à treize ans ét
002. 

de ces deux affaires, une homme et une FEMME : Alfred Sirven, ancien directeur
003. L
     001. Ult life of the average American man or WOMAN. Following a vigorous campaign of
004.
      002. In not standing up to them. An aggressive WOMAN wants a man to demand, not knuckle
005.
     003. 
all the way home. The ordinary man and WOMAN, however, saw little of the great
006.
     004. 

to the life of the middle-class man or WOMAN, dictating the methods of child re-
007.
008. 005. t exasperating to men, who expect every WOMAN to verify their preconceived notio
009.
     006. 
er, he muttered, to meditate on man (or WOMAN) than on God. David finished shavi
010.
      007. 🗆 dresser". "I guess it's children make a <u>WOMAN</u> old. A <u>man</u> gets old anyhow". After
011.
      008. 

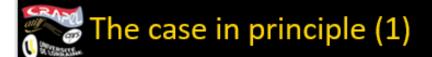
e. There was a very old man and a young WOMAN and a brood of children ranging fr
012.
     009. Trilted back. "So you're looking for a WOMAN who married a man who might have l
013.
014. [010. | ecial reason to believe that the man or WOMAN they sought had stayed only overni
015. \square 011. \square at? The door opened and three men and a WOMAN in a sari swept past him and down
016.
      012. 

a scout ship, partnered with a man or a WOMAN, whichever she chose, as the mobil
017.
      013. Deliver house. Now he saw that both the man and WOMAN were moving slowly and irregularly
      014. many affairs of the heart, a man and a WOMAN meet and something clicks. Somethi
019.
020. [015. ] in stormy weather I marry this man and WOMAN together. Let none but Him who rul
021. [
     016. Who rules the thunder Put this man and WOMAN asunder". Absolution for his lie?
022.
             fteen miles east of here. Eleven men, a WOMAN and a teen-age boy tramped over co
023.
             ." is a problem piece about a man and a WOMAN and the three "figures" that bothe
024.
     019. 

hem somehow. Unfortunately, the man and WOMAN were not made to appear very inter
025.
      020. 

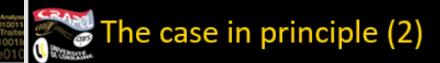
the rain came more heavily, and men and WOMEN in light summer clothes began to d
      021. 
on my experience to find as many men as WOMEN in church, and to hear almost ever
             he church school, missions, men's work, WOMEN's work, vouth program, social acti
```

Why *should* this learning method be effective?



Data-driven learning (DDL): using the tools and technique for teaching/learning/using a foreign/second language

- Greater effort better retention/deeper learning
- Strategy training
 learn to observe patterns in environment (exposure)
- Incidental learning
 "pick up" collocations while assembling word meanings
- Appeal to research-oriented clientele
 as opposed to ESL textbooks; PhD students, etc.



Data-driven learning (DDL): using the tools and techniques of CL for teaching/learning/using a foreign/second language

- A basis for CALL that goes beyond read/listen + MCQ with fixed answer
- Show not tell (constructivism)
 complex, dynamic, probabilistic, usage-based
- Autonomy own questions, life-long learning, multiple affordances
- Shapes up what students do anyway use Google for language research

But *learning* with corpora is not the topic of today's lesson

Rather *teaching* with corpora

Or, the design of courses conceived as corpora

And for this we need some DDL research tools

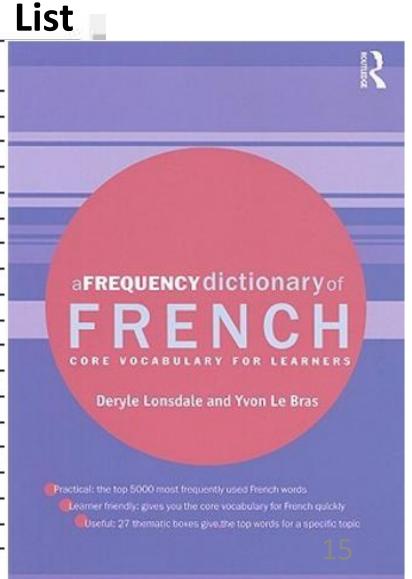
- Frequency lists based on both standard and course corpora
 - Frequency-list based tests
 - Frequency-list based text profiler
 - Lexical frequency profiler (LFP, Laufer & Nation)
- Some empirical findings to guide tool development
 - E.g. that reading comprehension depends on 95% know-word coverage
 - Laufer (1989 ... 2020)

The standard French corpus used on Lextutor

Corpus

	# de mots (approx.)	Genre
Oral		
	175,000	Conversations
	3,750,000	Hansard canadien
Y	3,020,000	Transcrits d'interviews
i i	1,000,000	Débats parlementaires de l'UE
	855,000	Appels téléphoniques
i i	470,000	Dialogues de théâtre
	2,230,000	Sous-titres de films
TOTAL	11,500,000	
Écrit		
	3,000,000	Agences de presse
	2,015,000	Articles de quotidiens
	4,734,000	Œuvres littéraires(fiction, non-fiction)
	434,000	Magazines de vulgarisation technique
	1,317,000	Bulletins de presse, manuels d'emploi
TOTAL	11,500,000	
AU TOTAL	23,000,000	

Tableau 1: Composition du corpus



FREQUENCY LIST

- 1. de 64006
- 2. la 31617
- 3. le 23365
- 4. et 20555
- 5. les 19532
- 6. des 19202
- 7. en 15185
- 8. du 14915
- 9. un 12781
- 10. a 12091
- 11. une 11168
- 12. est 10026
- 13. qui 9128
- 14. que 8796
- 15. dans 8124
- 16. pour 7600
- 17. par 7469
- 18. au 6991
- 19. pr 6428
- 20. sur 6224
- 21. pas 5382
- 22. ne 5303
- 23. plus 4987

Lemmatized and sorted by LEMMA frequency (not individual word freq.)

Aller 10,000
vais 3,000
vas 750
allé 500
etc.
TOTAL 14,250

... and divided into groups of 1,000 lemmas

etc ... \rightarrow 25k 2k 3k 4k https://lextutor.ca/vp/c lextutor.ca/vp/comp/fr 5 lextutor.ca/vp/c lextutor.ca/vp/com échéant écart ère âgé échelon échanger échéance âme à éclaircir échouer échantillon écarter âge économiquement éclairer éclat échange école économiser écouler écologique échapper économie écrasant économiste écran échec économique écrouler écraser écoute échelle élégant écouter écrit écrier écho écrire écriture édifice éloge éclater émeute éducation éditorial écu écrivain éminent également éditeur égyptien édition énergique élargir égard élémentaire égal énumérer électricité élaboration élément égalité épanouir électrique élan église élection épouvantable électronique élargissement élève élever équation élite élevé élaborer émission équivaloir émotion élimination. électeur énergie ériger énormément élu électoral époque étroitement émaner éliminer épargne équipe évasion émerger épouser élire équipement éventail éauilibrer émetteur éloigner établir éventualité émouvoir éauité émettre établissement abord équitable énergétique énorme état aboutissement énoncer étage épaule étranger abstenir éteindre épais époux étude accélération éternel épargner épreuve étudiant acceptation éthiaue épidémie éprouver étudier accessoire étoile épisode équilibre accompli événement étonnant épuiser équivalent accumulation éviter étroit équipage été acharner être équiper évacuer étape acheteur éventuellement accès étaler étendre

état major

aîné

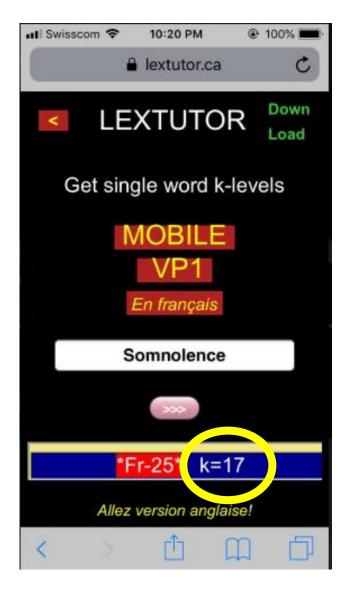
átonnar

accepter

acide

Such that every word has its 'serial number'





Reflecting the order they are typically found in the environment

+ the importance of their being learned

(depending on the corpus used, obviously)

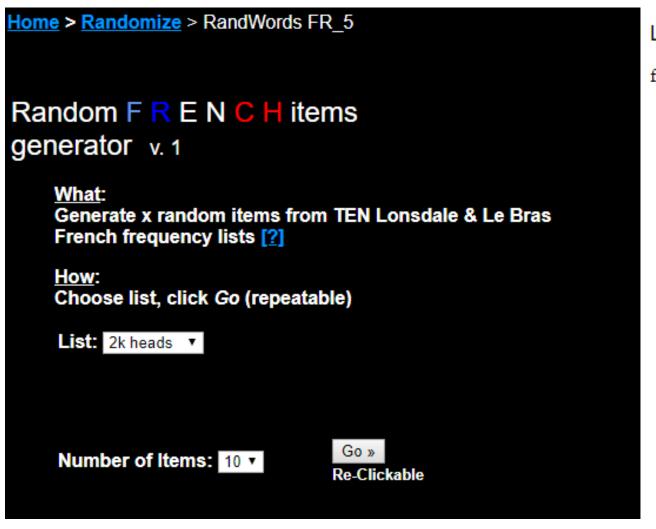
So, with corpus and numbered word-lists in hand... What can we do with this?

 Mainly, we have the means of providing learners with systematic vocabulary development

 Starting with systematic tests of vocabulary level in any language

With these lists and some tools anyone can make reliable level-based tests

ONE
 Format Oui-Non



List framework IS FR_5
fr 5 heads(02).txt

vide fabriquer pétrole revoir impression ignorer acteur entretien avocat côte

List framework IS FR_5

fr_5_heads(02).txt

poulard sconber radeurs cageait atérait

vide
fabriquer
pétrole
revoir
impression
ignorer
acteur
entretien
avocat
côte

USE RANDOM WORDS TO ...

OPTION ONE

Do whatever you like

OPTION TWO

Make novel Yes-No Checklist Tests in four easy steps (Why & How)

[1] From Generator on the left screen, get 25 random real words at a level, eliminate 5

[2] Click HERE to add 50% PNWs to your list (Length-appropriate random PNWs)

See all PNWs 'plausible non-words'

te which are PNWs before integrating!

[3] Click HERE to integrate

[4] Make an MS-Word document with three or four columns or use template

List framework IS FR 5

atérait

cageait

entretien

fabriquer

impression

ignorer

poulard

pétrole

radeurs

sconber

revoir

vide

avocat

côte

fr_5_heads(02).txt

USE RANDOM WORDS TO...

OPTION ONE

Do whatever you like

OPTION TWO

Make novel Yes-No Checklist Tests in four easy steps (Why & How)

[1] From Generator on the left screen, get 25 random real words at a level, eliminate 5

[2] Click HERE to add 50% PNWs to your list (Length-appropriate random PNWs)

See all PNWs 'plausible non-words'

Note which are PNWs before integrating!

[3] Click <u>HERE</u> to integrate

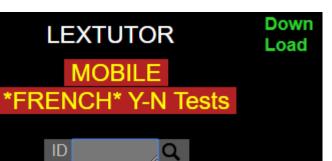
[4] wake an MS-Word document with three or four columns or use template

[5] Select, Copy, & Paste your words to

FRANCAIS 101 - 2k receptive vocab. Knowledge/diagnostic test PROF BATISTA / student _____ 21 MAR 2019

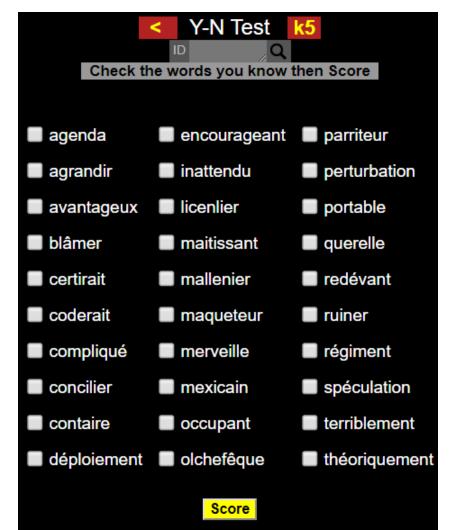
Put a check mark beside each word you know, however slightly. Some of the words are not real words, do not check those.

2k	3k	4k	5k
acteur	aistrope	abbey	aerial
atérait	axis	abolish	bourgeois
avocat	battalion	acklon	cambule
cageait	berrow	blunt	cramp
côte	beta	bypass	dole
entretien	binoculars	clap	dowrick
fabriquer	charitable	connery	eccentric
ignorer	commemorate	cradle	eldred
impression	dazzle	delete	fatigue
poulard	dispense	dissent	gummer
pétrole	eldred	duffin	haque
radeurs	elusive	elastic	hinge
revoir	forthcoming	immerse	kiley
sconber	galpin	junction	lavish
vide	inhale	lannery	leaflet
	knight	limp	millilitre



100 and 1,000-Lemma Sets x Freq.

K0-1	K1-2	K2-3	K3-4	K4-5	
C1	C1	C1	C1	C1	
C2	C2	C2	C2	C2	
C3	C3	C3	C3	C3	
C4	C4	C4	C4	C4	
C5	C5	C5	C5	C5	
C6	C6	C6	C6	C6	
C7	C7	C7	C7	C7	
C8	C8	C8	C8	C8	
C9	C9	C9	C9	C9	
C10	C10	C10	C10	C10	
K1	K2	K3	K4	K5	
List Info					
"k1c1" = 1,000-1,100 zone					
Allez version anglaise!					
<www.lextutor.ca< td=""></www.lextutor.ca<>					





Or something more standard, similar to Nation's VST

• Exemple (Le "TTV", RCLV, 2016)

A New Receptive Vocabulary Size Test for French

Box 4. A noun cluster from 5K frequency section of the TTV 1. brouillard 2. coïncidence 3. farce _____ une histoire qui fait rire 4. instituteur ____ ce qui empêche de voir loin 5. pneu ____ un professionnel de l'éducation 6. soumission

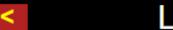
A New Receptive Vocabulary Size Test for French

Roselene Batista and Marlise Horst

Abstract: Researchers have developed several tests of receptive vocabulary knowledge suitable for use with learners of English, but options are few for learners of French. This situation motivated the authors to create a new vocabulary size measure for French, the Test de la taille du vocabulaire (TTV). The measure is closely modelled on Nation's (1983) Vocabulary Levels Test (VLT) and follows the guidelines written by Schmitt, Schmitt, and Clapham (2001). Initially, a pilot version was trialled with 63 participants; then an improved version was administered to 175 participants at four proficiency levels. Results attest to the TTV's validity: mean scores across the four frequency sections decreased as the tested words became less frequent, and more proficient learner groups outperformed less proficient groups. The TTV in its current form is intended to be of practical use to teachers and learners, but it is also expected to evolve; ideas for future improvements are discussed.

Keywords: frequency, French L2 vocabulary, vocabulary size, assessment

Résumé: Des chercheurs ont développé plusieurs tests de vocabulaire réceptif pour les apprenants d'anglais, mais les options pour les apprenants de français ne sont pas nombreuses. Ce scénario a motivé les auteurs à créer un nouvel outil qui mesure la taille du vocabulaire en français, le Test de la taille du vocabulaire (TTV). Cet outil repose sur le modèle du Vocabulary Levels Test, conçu par Nation (1983), et suit les directives proposées par Schmitt, Schmitt et Clapham (2001). Premièrement, une version pilote a été testée auprès de 63 participants, ensuite une version améliorée a été complétée par 175 participants de quatre niveaux de compétence distincts. Les résultats confirment la validité du test: les moyennes obtenues par les participants à chacune des qua-



LEXTUTOR

Down Load

MOBILE LEVELS TESTS

Score-Recording Tests

Back Button to return to this page

VLT

Vocab Levels Test

VLT_{v.2}

VST

Vocab Size Test

TTV

Test de la taille du vocab

<www.lextutor.ca>



	SC	C	R	E	
		a			
	pr	ac	tic	e	_
2k	sc	or	e		
	0 (
Pro					
2	3	4	5	6	7
3k	sc	or	e		
0/3					
Pro	0002				
1	2	3	4	5	6
5k 0/3					
Pro	ob.	S	ets		
1	2	3	4	5	6
10	ks	co	re		
0/3	0 (05	6)		
Pro	ob.	5	ets		
1	2	3	4	5	6
	ET		RE	Œ	R
	1>	1			
	3>4				
	5>(3			
	7>1				
	9>	10			

Typical findings from either test are that ~

- Intermediate learners know ≈ 1,000-1,500 lemmas in their L2
- High school graduates know ≈ 2,000-3,000 lemmas

- But is this good or bad?
 - Sufficient or insufficient?
 - To do what?

To answer this, we also need a similar *test* for texts

On the same metric as the tests for learners

- A.k.a. a text profile
 - How many words at each k-level does this text contain?
 - 'Know' replaced by 'contain'

Text profiling by frequency level

 Also known as LFP and Vocabprofil(e) – "VP"

- For example, this text (893 words) from a Montreal news story
 - - "L'abandon scolaire"

Sus à l'abandon scolaire! mais je ne pense pas qu'on 200 écoles second e mettre Profile summary partageront 125 m e ce que le is.» Est-il trop tôt cinq ans K-1 816 85.2 Marie-Andrée Cho sque K-2 66 92.1 Le mardi 14 mai 2 pilote n'a pas K-3 36 95.9 Pour s'attaquer au sa première K-4 11 97.0 qu'est l'abandon s ucun résultat 97.6 K-5 6 réduire de 10 % er disponible? K-6 97.7 ministère de l'Édu hmes pas encore 97.9 K-7 (MEQ) propulse 12 lire avec précision en cing ans dans u K-8 98.1 impacts du de 200 écoles seco mais nous en K-10 98.5 Agir autrement, c' pour dire que 98.7 K-11 de l'opération land portante, et le K-12 98.9 le ministre de l'Édi ble», notait 99.1 K-13 Sylvain Simard, ins on, sous-ministre K-14 99.2 droit d'un projet-p seignement K-17 99.3 même nom, qui n' condaire, K-18 99.4 livré ses premiers le suivi de ces L'intervention ne b K-20 99.5 serré. Et n de fonds? «Ce toutefois pas de la K-24 99.6 enveloppe qu'ont ement une ≈100 OFF 3 écoles secondaires gent», ajoute 🎪

l'expérimentation: chacune a

Bisaillon. «Je souhaiterais que

Colour-coded k-levels

Profile summary K, #tokens, cumul%					
K-1	816	85.2			
K-2	66	92.1			
K-3	36	95.9			
K-4	11	97.0			
K-5	6	97.6			
K-6	1	97.7			
K-7	2	97.9			
K-8	2	98.1			
K-10	4	98.5			
K-11	2	98.7			
K-12	2	98.9			
K-13	2	99.1			
K-14	1	99.2			
K-17	1	99.3			
K-18	1	99.4			
K-20	1	99.5			
K-24	1	99.6			
OFF	3	≈100			



sus à le abandon scolaire chiffre écoles secondaires se partageront chiffre millions en cinq ans marie andrée chouinard le mardi chiffre mai chiffre pour se attaquer au fléau que est le abandon scolaire et le réduire de chiffre en dix ans le ministère de le éducation meg propulse chiffre millions en cing ans dans un concentré de chiffre écoles secondaires agir autrement ce est le nom de le opération lancée hier par le ministre de le éducation sylvain simard inspirée tout droit de un projet pilote du même nom qui ne a pas encore livré ses premiers résultats le intervention ne bénéficiera toutefois pas de la généreuse enveloppe que ont reçue les six écoles secondaires ciblées par le expérimentation chacune a reçu chiffre chiffre million pour des mesures échelonnées sur trois ans tandis que les chiffre écoles ciblées par québec se disputeront chiffre millions par an distribués selon le bon vouloir de leurs commissions scolaires respectives pour permettre cette annonce la directrice de le école secondaire Édouard montpetit lucie lalande ouvrait sa porte au ministre sylvain simard hier depuis presque un an son école est sous la lorgnette du meq parce que elle fait partie du projet agir autrement coups de centaines de milliers2de dollars les résultats sont déjà palpables affirme elle

Compare previous profile to this one ~

 Hier Daisy s'était levée tôt ce matin de printemps. Elle travaillait sur une affaire dans la ville voisine. Elle arriva à son bureau à huit heures avec à la main un sac en papier contenant des petits pains. Elle mourrait d'envie d'une tasse de café.

Note – Demo only; short texts = unreliable VP

Profile summary K, #tokens, cumul%					
K-1	42	87.5			
K-2	3	93.7			
K-3	2	97.9			
K-6	1	100.0			
OFF	0	≈100			

hier daisy se était levée tôt ce matin de printemps elle travaillait sur une affaire dans la ville voisine elle arriva à son bureau à huit heures avec à la main un sac en papier contenant des petits pains elle mourrait de envie de une tasse de café

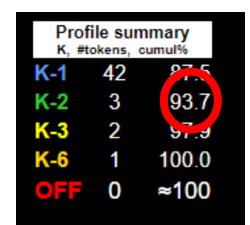
Profiles are somewhat different!......

But without these instruments, would we have noticed how much?

To read Daisy with 95% word knowledge, 2000 word families should be know

To read Abandon Scolaire with 95% word knowledge, 3000 word families should be know

Daisy le matin



L'abandon scolaire

Profile summary K, #tokens, cumul%				
K-1	816	85.2		
K-2	66	02.1		
K-3	36	95.9		
K-4	11	57.0		
K-5	6	97.6		
K-6	1	97.7		
K-7	2	97.9		
K-8	2	98.1		
K-10	4	98.5		
K-11	2	98.7		
K-12	2	98.9		
K-13	2	99.1		
K-14	1	99.2		
K-17	1	99.3		
K-18	1	99.4		
K-20	1	99.5		
K-24	1	99.6		
OFF	3	≈100 31		

These tests and measures let us compare/contrast profiles of learners & profiles of texts

And determine the quality of the match

- Empirical research has shown
 - Ss need to know 95% of the words on a page to read with basic comprehension
 - > 95% of words known = fluent reading for comprehension
 - **85-95%** = intensive reading
 - Many look-ups, comprehension gaps, unresolved passages
 - < 85% = unuseful experience
 - (zone where much of school reading lies, in L1 but especially L2)

With these tools we can begin Data-Driven Course/Curriculum design (DDCD)

Start with the raw materials of a course

 E.g., A "reading course" with a set of texts drawn as "interesting" from Internet news

Step 1 – get all course materials into digital form as a corpus

- In two formats:
 - Separate chapters, texts or units (ideally 'zipped')
 - 2. One large single file
 - What will this enable us to do?

What would course-as-corpus let us do?

Find out

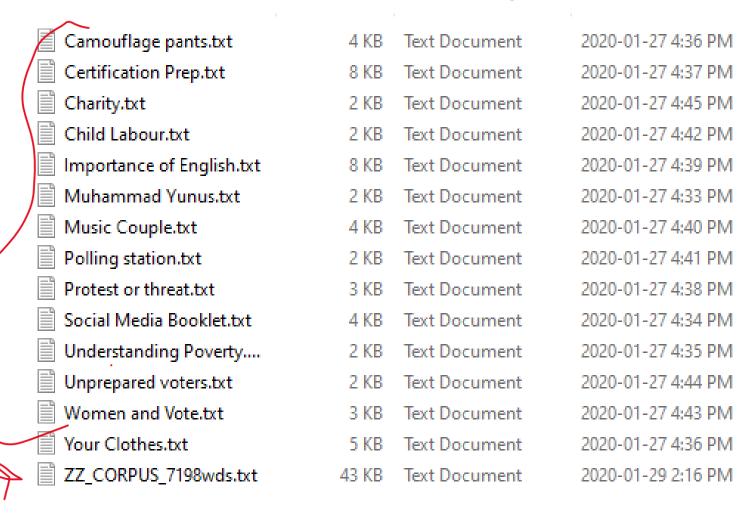
- 1. Whether our course level is anywhere near our learners' level
- 2. If not, what would bring it nearer?
 - 1. Learners up, texts down?
- 3. What if anything can legitimately be tested from this course?

^{*} Level in this case refers to lexis but could be other aspects

Here is the course-corpus described in one course I re-designed

ESL Adult Reading, 14 units,

- >7000 words,
- in-house course,
- intermediate adults
- high stakes
- → Simple 'flat' corpus
 - 2 formats
 - Separate files
 - A combined file



Now a new course For FL2 learners in Sweden

Le Petit Nicolas

(Sempé and Goscinny, c. 1959)



- 18 children's stories
- Similar lengths
- No passé simple/historique
- 28,000 words (≈ 1500 per story)
- a plausible French reading course for 9-year-olds

Name	Date modified	Size	Туре
01_souvenir_cherir.txt	2022-09-01 5:29 PM	8 KB	TXT File
2 02_cow-boys.txt	2022-09-01 5:30 PM	8 KB	TXT File
3_le_bouillon.txt	2022-09-01 5:30 PM	8 KB	TXT File
04_le_football.txt	2022-09-01 5:31 PM	8 KB	TXT File
05_l'inspecteur.txt	2022-09-01 5:31 PM	8 KB	TXT File
	2022-09-01 5:32 PM	8 KB	TXT File
07_djodjo.txt	2022-09-01 5:32 PM	8 KB	TXT File
08_chouette_bouquet.txt	2022-09-01 5:32 PM	8 KB	TXT File
09_les_carnets.txt	2022-09-01 5:32 PM	8 KB	TXT File
10_louisette.txt	2022-09-01 5:33 PM	14 KB	TXT File
11_je_fume.txt	2022-09-02 11:39 AM	8 KB	TXT File
12_le_petit_poucet.txt	2022-09-01 5:34 PM	8 KB	TXT File
13_le_velo.txt	2022-09-02 11:39 AM	8 KB	TXT File
14_suis_malade.txt	2022-09-01 5:34 PM	8 KB	TXT File
15_bien_rigole.txt	2022-09-01 5:34 PM	7 KB	TXT File
16_frequent_agnan.txt	2022-09-01 5:34 PM	9 KB	TXT File
17_pas_le_soleil.txt	2022-09-02 10:37 AM	8 KB	TXT File
18_quitte_la_maison.txt	2022-09-02 12:01 PM	8 KB	TXT File
Pet_Nic_18.zip	2022-09-02 12:02 PM	61 KB	ZIP File
PN_original_30k.txt	2022-09-05 9:53 AM	136 KB	TXT File

Petit Nicolas (tout) - lexical frequency profile

Home > VocabProfilers > VP-Compleat Input > Output ('Back' preserves inputs) FRAMEWORK IS fr_5

Freq. Level	Flemmas (%)	Types (%)	Tokens (%)	Cumul. token (%)
K-1:	2451 (72.4)	2651 (65.98)	23531 <u>(86.3)</u>	86.3
K-2:	248 (7.3)	454 (11.30)	1420 <u>(5.2)</u>	91.5
K-3:	157 (4.6)	256 (6.37)	660 <u>(2.4)</u>	93.9
K-4:	110 (3.2)	152 (3.78)	354 <u>(1.3</u>)	95.2
	Co	overage 95	[?]	
K-5:	83 (2.5)	111 (2.76)	263 <u>(1.0</u>)	96.2
K-6:	92 (2.7)	122 (3.04)	274 <u>(1.0)</u>	97.2
K-7:	36 (1.1)	52 (1.29)	150 <u>(0.5)</u>	97.7
K-8:	33 (1.0)	39 (0.97)	93 (<u>0.3)</u>	98.0
	Co	verage 98		
K-9:	26 (0.8)	29 (0.72)	72 (<u>0.3)</u>	98.3

Result

95% lexical coverage is reached only when 4,000 lemmas are known

Knowing 1,000=86.3% coverage

Knowing 2,000=91.5% coverage

How many French words do Swedish 9-yr-old learners know? Probably < 1,000

So these texts probably correspond to learners' interests but...

- Will be challenging for some/most of them
 - (This could only be established through testing)
 - So here is a prima facie test case for Data-Driven Course Deign

How can DDL help make this text accessible?

Simplest, DDL can just give us the specific 2k-3k words to teach

- Whether Pre-Teach or emphasize and invite discussion when encountered
- Teacher would choose based on frequency, cognate-ness
- For the whole book or specific story/stories
 - (Next slide is for whole book)

French v.5 -K2 Flemmas: [flems 248 : types 454 : tokens 1420] VP-negative: fr_5-2 absolument_[3] accompagner_[3] acteur_[1] agent_[5] anglais_[1] animal_[3] apercevoir_[4] appareil_[2] approcher_[7] arbre_[7] arrière_[1] assister_[1] attacher_[3] attaquer_[2] avion_[15] baisser_[1] balle_[15] French v.5 -K3 Flemmas: [flems 157 : types 256 : tokens 660] ca cie VP-negative: fr 5-3 CO agiter [7] allumer [1] amuser [32] animer [1] anniversaire [3] arracher [1] arranger [6] bouteille [8] boîte_[6] brillant_[2] briller_[4] bête_[7] cadeau_[5] calmer_[1] camarade_[13] capitaine_[7] chaise_[3] French v.5 -K4 Flemmas: [flems 110 : types 152 : tokens 354] corrig esi dent fin VP-negative: fr_5-4 douce ga embri accueil [1] affreux [1] alcool [2] amateur [1] applaudir [1] attraper [1] auto [21] avaler [1] inc avertissement_[1] bain_[6] ballon_[8] banc_[13] blague_[4] bouton_[2] bravo_[2] brutal_[3] bâton_[1] cahier [1] canon [1] caresser [2] chant [1] charmant [1] chasseur [3] chemise [10] chimie [1] château [1] ma chèque_[1] coincer_[1] col_[1] combiner_[2] consoler_[2] corde_[2] correctement_[1] costume_[5] cou_[1] crâne_[1] dicter_[1] déboucher_[1] débrouiller_[2] déchirer_[3] dégât_[1] délivrer_[1] dépêcher_[1] pa morci désespérer [2] désordre [1] détacher [1] essuyer [6] fantôme [2] ficher [1] fréquent [2] fréquenter [2] fusil_[1] fêter_[1] gendarme_[3] gorge_[1] goûter_[9] guérir_[2] habiller_[14] herbe_[3] hurler_[2] jouet_[12] joyeux [1] lampe [2] laver [3] liquide [2] lunette [24] maigre [1] mairie [1] manche [2] miroir [1] mordre [8] méfier [4] métal [1] nerveux [7] nettoyer [1] os [1] pantalon [4] patte [2] peindre [1] pendre [2] rer photographe_[15] photographie_[1] piquer_[1] planche_[3] plume_[3] poing_[19] poitrine_[2] pot_[2] proclamer_[1] punir_[6] pâle_[1] queue_[4] rage_[1] rattraper_[4] rembourser_[1] ressort_[1] sourd_[2] se stupide_[1] sympathie_[1] tache_[3] tante_[5] taper_[15] thé_[3] transparent_[1] vitre_[1] voleur_[5] élan_[1] élémentaire [1] SU traduire [

voisin_[4] voier_[2] voyage_[1] zero_[1] ecarter_[1] ecnapper_[1] eciater_[1] eioigner_[1] etonner_[12]

A way to find the most important words Ss must know across a set of texts is with the DDL tool RANGE

- Gives the frequency of each word or lemma in a set of texts
 - Plus the number of texts each features in
- Can boil the lexis down to the frequent PLUS recurring items

Home > Range > Range for Texts NEW 2020: APRIL - zip uploads; NOV - file freqs

Range for Texts v.5.3

Estimated Capacity 8 novels x 125k wds = 1,000,000 wds(/40 secs)

Upload up to 25 text files (chapters of a book, works of an author, frequency lists) and then: (1) find the frequency of each word in the c range of each across texts (e.g., is in 6 out of 9 texts), (3) build a Frequency x Range word list for the collection [?]

Home> Range Input >[Use «Back to preserve inputs]> Range for Texts v.5.1 - Output VP-Frame See I bottom for filenames, stoplists, offlist, profile, output filtering (freq x range), & Excel-copiable version of TITLE: Pet Nic_18.zip Post-Analysis=> VP (token coverage) Range Profile Loading... Click headings to sort ◆ (In the table, T1(25) means text T1 has 25 occurrences of the word, etc) 000. Fams Freq Range VP T1 T2 T3

032.	papa	172	16	2	T1(2)	T2(16)	
036.	maman	154	13	2		T2(5)	
041.	maître	118	12	1	T1(27)		T3(8)
061.	crier	80	16	2	T1(6)	T2(4)	T3(10)
075.	répondre	50	17	1	T1(3)	T2(2)	T3(4)
082.	école	46	15	1	T1(2)		T3(2)
084.	pleurer	44	17	2	T1(2)	T2(1)	T3(6)
088.	fleur	42	2	2			
092.	vélo	40	3	5			
093.	bouillon	39	3	11			T3(32)
095.	rex	39	1	10			
100.	chouette	36	17	7		T2(5)	T3(1)
103.	copain	35	18	3	T1(3)	T2(2)	T3(1)
110.	rigoler	33	13	5	T1(1)	T2(4)	T3(1)
111.	drôlement	32	15	7	T1(3)	T2(1)	T3(2)
115.	inspecteur	30	2	3			
116.	tête	30	12	1	T1(3)	T2(5)	
118.	amuser	29	13	3		T2(4)	
122.	courir	28	13	2	T1(1)	T2(2)	
129.	idée	26	13	1	T1(1)	T2(2)	T3(1)

43

RANGE EXTRACT || FREQUENCY => 7 and RANGE => 10 || 24 i

Offlist

75 OFF-LIST TYPES WITH NO FAMILY HEADER

ac agnan alceste aon arcachon australie blédurt bordenave bouil boxing bril cyrille dgeorges dinguedingue dissipez djoachim djodjo dubon dég elvire eude

OUTPUT FILTER





Extract data by Frequency >= 7 v and Range >= 10v

Format: [• As above • Words only] Go »

Build Frequency- and Range-based word lists; check corpus

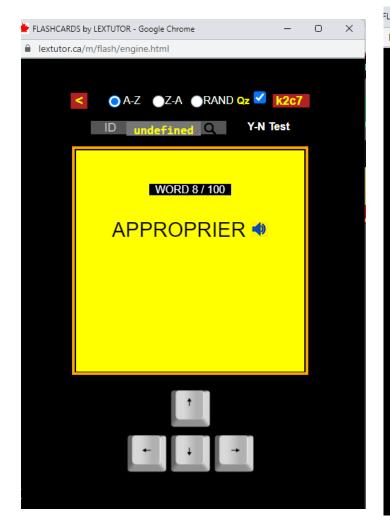
These are the key words across these stories –

List can be expanded or shortened by changing the parameters

				_
ı	-AM	FREQ	RANGE	V
035.	maman	154	13	2
040.	maître	118	12	1
060.	crier	80	16	2
074.	répondre	50	17	1
081.	école	46	15	1
083.	pleurer	44	17	2
099.	chouette	36	17	7
102.	copain	35	18	3
109.	rigoler	33	13	5
110.	drôlement	32	15	7
115.	tête	30	12	1
117.	amuser	29	13	3
121.	courir	28	13	2
128.	idée	26	13	1
135.	chocolat	24	12	5
151.	arrêter	21	11	1
164.	rouge	19	10	2
165.	terrible	19	11	2
168.	embêter	18	11	6
175.	gentil	17	10	3
187.	sûr	16	13	1
196.	côté	15	10	1
197.	drôle	15	10	2
204.	taper	15	11	4

Once the key words' identities are known...

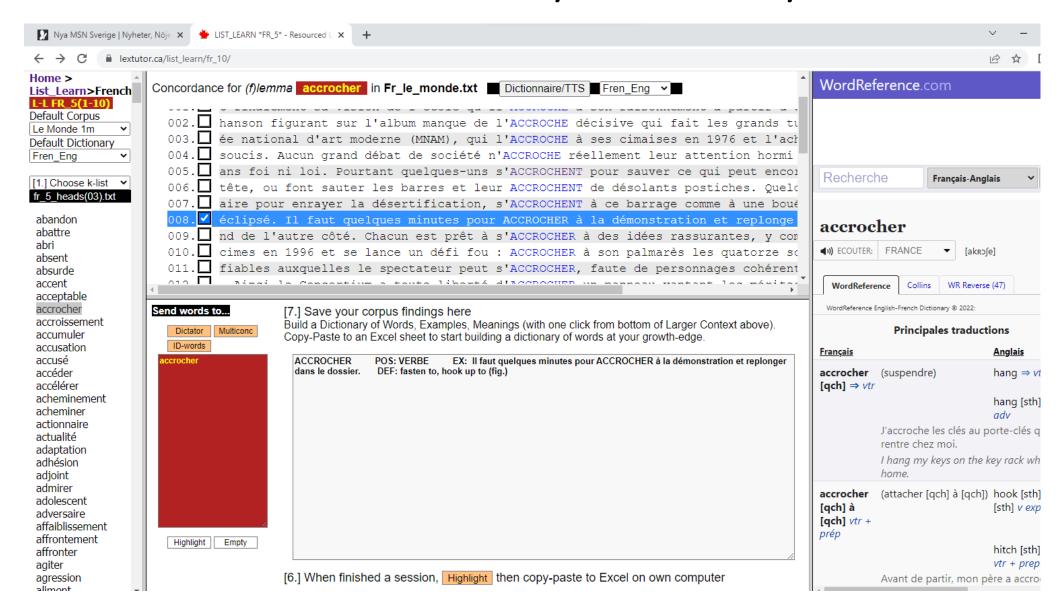
Flash Cards are back in style with shiny new credentials







Personally I would have Ss go through lists and check out *all the* words at the next level that they don't already know



More student-centered is to set up a Group Lex database & Ss choose their own words to work on

- And share with others
 - Including by mobile phone
- Leading to a complete lexicon of the whole book with quizzes

Accueil >Group Lex		ous		MOTS: Op st 10 2 nd 10		uiz>> I Imprimer Interactive 4^{th} 10 5^{th} 10 6^{th} 10 7^{th} 10 8^{th} 10	_	10 10 th 10 20 au hasard	Aucun	Création du q	uiz
v.8.3		#	Qz	NOUVEAU MOT ◆	VP ♦	EXEMPLE	CLAS DE MO		GROUPE \$	ETUDIA	DATE/H SOUMIS
Beta French interface		1	/		2	« La mi-temps de quoi? a demandé Alceste. Je viens de m'apercevoir que nous n'avons	Verhe	skymta, glimt	Det Nic /	tom_ad	2022.09.06
OS=Win BROWSER=Chr	\			apercevoir	2	pas de ballon, je l'ai oublié à la maison! »	Verbe	skymica, giime	rec_mc_s	r tom_au	8.30
V.Mobile		2		musique	1	J'écoute de la musique.	Nom	L'art de la combinaison de sons	Arts	LK1234	2014.12.16 9:31
pad/phone v. English		3		couple	2	Il vit en couple avec Marie.	Nom	deux personnes liées par des sentiments	Economic	E LK1234	2014.12.16 11:38
COMPUTER Voir les		4		<mark>couple</mark>	2	Seulement un couple sur trois désormais se marie à l'église.	Nom	deux personnes liées par un sentiment	Economic	123French	2014.12.16 11:32
54 mots Ajout. mot		5		<mark>pris</mark>	1	Elle a pris sa douche antérieurement	Verbe	took (prendre)	Arts	00christi	2014.09.08 9.35
Aj. étudiant Modif. mot		6		pomme	4	La pomme de pin n'est pas un vrai pomme	Nom	Apple	Arts	tom_ad	2014.09.08 9.33
Besoin de Gp		7		montre	1	Que dit ta montre?	Nom	(wrist) watch - so 'What time is it?'	Arts	tom_ad	2014.09.08 9.28
X-tract-liste		8		nombreux	1	On est demandé de venir 'nombreux' à leur évenement	Nom	in large numbers	Arts	tom_ad	2014.09.08 9.24
Recherche?		9		<mark>salut</mark>	2	Dis Salut	Nom	Hello	Arts	tom_ad	2014.09.08 22:34
Prononcer la sélection		10		sourire	1	Donne-moi ta sourire	Nom	A smile	Arts	tom_ad	2014.09.08 22.23

Accueil >Group Lex

v.8.3

Beta French interface

OS=Win BROWSER=Chr

> V.Mobile pad/phone

v. English

COMPUTER Voir les 53 mots

Ajout. mot

Aj. étudiant

Modif. mot

Besoin de Gp Lx?

X-tract-liste

Recherche?

Prononcer la sélection

Ceci sera votre entrée :

#	Qz	NOUVEAU MOT	VP	EXEMPLE	CLASSE DE MOT	DÉFINITION	GROUPE	ÉTUDIANT	DATE/HEURE
0		apercevoir		« La mi-temps de quoi? a demandé Alceste. Je viens de m' apercevoir que nous n'avons pas de ballon, je l'ai oublié à la maison! »	V	skymta, glimt	Pet_Nicolas_4	tom_admin	2022.09.06 04.23

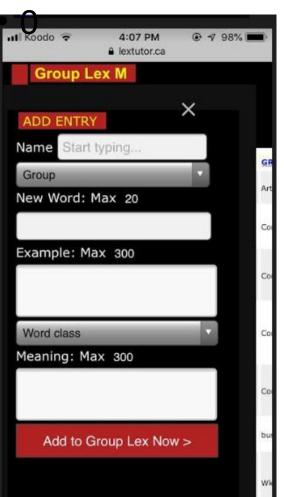
< Retour

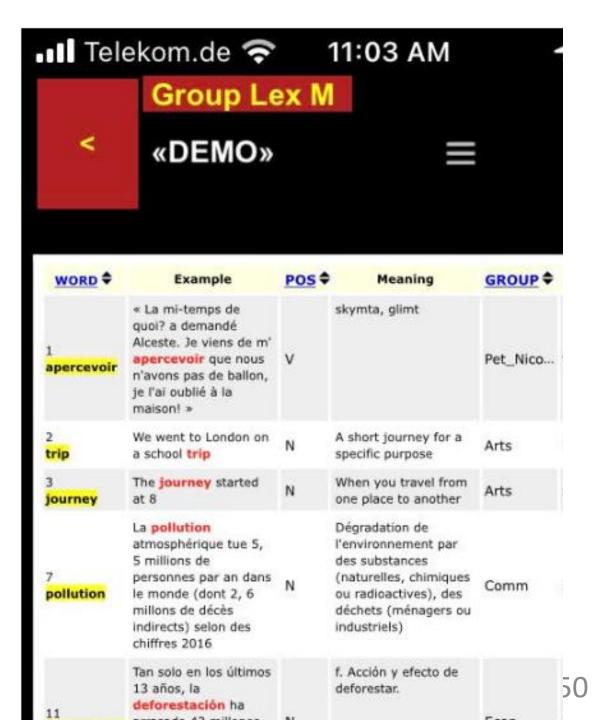
Sécurité : Taper les chiffres pour trois, deux, un



Valider >

XXX





QUIZ 1 - contexts from WHOLE GP LEX 08 Feb 20, 16:49 Quiz 2 -new contexts - after 100% of History >> *Check* Questions: 10 Correct: 0 Tries: 0 Percent: 0 Practice In-Class 🎩 WORD NEW **EXAMPLE** DEF WORD CLASS A division or contrast I Rousseau's fundamental idea was the are or are represented of nature and culture. The word 'flight' has two different NEW WORD ambitious **EXAMPLE** WORD CLASS plane journey, and the act of running av ambitious He This was an ADJ person project. requiring a great de He was an person eagerly desirous of asteroid He was an person wealth, a specific of House prices have been dichotomy Aerial of the enemy position showed they were ADJ Information gatheri He This was an pers ready to attack. meaning Gp Lex also works on sends a beacon signal to the ground mobile indicating how urgent it is to track the spacecraft for a vehicle used for t reduced functionality N telemetry. Aerial of the enemy reconnaissance ey Gp Lex also works on though with reduced were ready to attack. Easily moveable functionality spacecraft sends a beacon The cientists are hoping an will clue them in about Ν one of many large indicating how urgent it is to track the sistationary early life on earth. ADJ House prices have been for months. not moving, or not A division or contra Rousseau's fundamental idea was the _____ of nature or are represented and culture. different. The word 'flight' has two different of somet s: a plane journey, and the act of running away. represents 字 We are collecting English one after another

MORE PRINCIPLED HOWEVER — is to use DDL tools to sequence the 18 stories from basic to more complex vocab

- To facilitate the conditions whereby learning by inference can take place
 - IE, > 90% knowledge of words in the context
 - Or as many as possible

These texts...

- Written for NS French kids, these are unlikely to be sequenced at all
- To get a sequence, we return to VP
- And analyse the texts individually, recording each of the 18 profiles
- Then with Excel we sort by the percentage of 1k items
 - Should give us a sequence of texts from most-basic to least-basic lexis
 - At least, this is what happened in English
 - Same in French?

K-1:	305 (73.3)	354 (72.69)	1309 <u>(86.2)</u>	86.2
K-2:	38 (9.1)	48 (9.86)	75 <u>(4.9)</u>	91.1
K-3:	18 (4.3)	21 (4.31)	28 <u>(1.8)</u>	92.9
K-4:	21 (5.0)	23 (4.72)	45 <u>(3.0)</u>	95.9
	С	overage 95	[7]	
K-5:	9 (2.2)	10 (2.05)	13 <u>(0.9</u>)	96.8
K-6:	4 (1.0)	4 (0.82)	4 <u>(0.3</u>)	97.1
K-7:	3 (0.7)	3 (0.62)	5 <u>(0.3</u>)	97.4
K-8:	2 (0.5)	2 (0.41)	2 <u>(0.1</u>)	97.5
K-9:	3 (0.7)	3 (0.62)	5 <u>(0.3</u>)	97.8
	Co	overage 98		
K-10:	1 (0.2)	1 (0.21)	3 (0.2)	98.0
K-11:	4 (1.0)	4 (0.82)	5 (0.3)	98.3
K-12:	2 (0.5)	2 (0.41)	5 (0.3)	98.6
K-13:	1 (0.2)	1 (0.21)	1 (0.1)	98.7
K-14:	2 (0.5)	2 (0.41)	8 <u>(0.5)</u>	99.2
K-15:	_ ()	_ (/	- <u>x</u> x	
K-16:				
K-17:				
K-18:				
K-19:				
K-20:	1 (0.2)	1 (0.21)	1 <u>(0.1)</u>	99.3
K-21:	. (5.2)	. (3.21)	- <u>X32.7</u>	55.5
K-22 :	1 (0.2)	1 (0.21)	1 (0.1)	99.4
K-23:	. (5.2)	. (0.21)	- <u>(3.7)</u>	33.1
K-24:				
K-25:	1 (0.2)	1 (0.21)	1 (0.1)	99.5
Off-List:	??	6 (1.23)	7 (<u>0.46</u>)	99.96
		J (1.25)	(0.10)	
Total (unrounded)	416+?	487 (100)	1518 (100)	≈100.00

VVOIUS III IGAL (IUNGIIS).	1310
Different words (types):	487
Type-token ratio (TTR):	0.32
Tokens per type:	3.12
Pertaining to onlist only	
Tokens:	1511
Types:	481
Flemmas:	416
Tokens per Flemma:	3.63
Flemma/token ratio (FTR):	0.28
Types per Flemma:	1.16
Singletons ratio	
Fams(n=1)[121] / total [416]	0.29

aste to Excel neet to record comparative coverages. To compare several texts, track learner writing over time, sequence texts by lexis (TIP, sort spreadsheet by K-1% to get more basic lexis first), etc. Note the FL7, etc.)

So 18 profiles assembled ready for sorting by Column 'C' - Order texts by the proportion of 1k

A	В	С	D	Е	F	G	Н	1	J	K	L	М	N	0	Р	Q	R	S	Т	U
LE PETIT NICOLAS																				
AS PUBLISHED																				
1 SOUVENIR CHERIR	K-1	86.2		K-2	91.1		K-3	92.9		K-4	<u>95.9</u>		K-5	96.8		K-6	97.1		K-7	97.4
2 COWBOYS	K-1	83.0		K-2	89.6		K-3	91.1		K-4	93.0		K-5	94.3		K-6	<u>95.0</u>		K-7	96.4
3 BOUILLON	K-1	88.2		K-2	92.4		K-3	94.2		K-4	94.9		K-5	<u>95.7</u>		K-6	96.5		K-7	96.8
4 FOOTBALL	K-1	86.4		K-2	92.6		K-3	<u>95.4</u>		K-4	97.1		K-5	97.9		K-6	99.2		K-7	99.3
5 INSPECTEUR	K-1	87.1		K-2	91.1		K-3	94.4		K-4	<u>96.0</u>		K-5	96.5		K-6	97.8		K-7	<u>98.4</u>
6 REX	K-1	82.6		K-2	90.6		K-3	92.6		K-4	93.7		K-5	93.9		K-6	94.6		K-7	<u>95.3</u>
7 DJODJO	K-1	87.0		K-2	91.0		K-3	94.2		K-4	<u>95.1</u>		K-5	95.6		K-6	96.0		K-7	96.7
0 8 CHOUETTE_BOUQUET	K-1	86.4		K-2	92.3		K-3	94.8		K-4	<u>95.7</u>		K-5	96.6		K-6	98.2		K-7	98.8
1 9 CARNET	K-1	88.6		K-2	93.8		K-3	<u>95.1</u>		K-4	95.5		K-5	97.6		K-6	<u>98.4</u>		K-7	98.7
2 10 LOUISETTE	K-1	87.5		K-2	93.4		K-3	<u>95.3</u>		K-4	96.3		K-5	97.0		K-6	97.7		K-7	<u>98.1</u>
3 11 JE_FUME	K-1	85.5		K-2	89.5		K-3	93.1		K-4	93.8		K-5	94.8		K-6	<u>97.0</u>		K-7	97.7
4 12 PETIT_POUCET	K-1	86.7		K-2	90.2		K-3	93.7		K-4	94.9		K-5	<u>95.3</u>		K-6	96.6		K-7	96.9
5 13 VELO	K-1	84.7		K-2	92.1		K-3	93.7		K-4	94.1		K-5	96.8		K-6	97.6		K-7	97.9
6 14 MALADE	K-1	85.5		K-2	90.4		K-3	93.8		K-4	<u>95.0</u>		K-5	97.1		K-6	<u>98.0</u>		K-7	98.6
7 15 BIEN_RIGOLE	K-1	88.2		K-2	92.3		K-3	94.4		K-4	94.9		K-5	<u>95.9</u>		K-6	96.4		K-7	96.9
8 16 FREQUENTE_AGNEN	K-1	85.9		K-2	90.7		K-3	93.1		K-4	<u>95.6</u>		K-5	96.6		K-6	97.4		K-7	<u>98.1</u>
9 17 PAS_DE_SOLEIL	K-1	86.4		K-2	90.8		K-3	92.8		K-4	93.9		K-5	94.7		K-6	<u>96.5</u>		K-7	97.1
0 18 QUITTE_MAISON	K-1	88.2		K-2	92.8		K-3	94.8		K-4	96.3		K-5	97.1		K-6	97.9		K-7	<u>98.4</u>
1																				

Giving us this...

Hmm, this is less useful than expected

1k-Range is only 6 percentage points over 18 stories – is it worth the trouble? Or do we need to know something about French to make this work

		• .		-	-	6				12		14	N.	0		0	-		_		
А	В	С	D	E	F	G	Н	1	J	K	L	M	N	0	Р	Q	R	S	Т	U	
LE PETIT NICOLAS SORTED BY 1k		1																			
4 9 CARNET	K-1	88.6		K-2	93.8		K-3	<u>95.1</u>		K-4	95.5		K-5	97.6		K-6	<u>98.4</u>		K-7	98.7	
5 3 BOUILLON	K-1	88.2		K-2	92.4		K-3	94.2		K-4	94.9		K-5	<u>95.7</u>		K-6	96.5		K-7	96.8	
6 15 BIEN_RIGOLE	K-1	88.2		K-2	92.3		K-3	94.4		K-4	94.9		K-5	<u>95.9</u>		K-6	96.4		K-7	96.9	
7 18 QUITTE_MAISON	K-1	88.2		K-2	92.8		K-3	94.8		K-4	<u>96.3</u>		K-5	97.1		K-6	97.9		K-7	<u>98.4</u>	
8 10 LOUISETTE	K-1	87.5		K-2	93.4		K-3	<u>95.3</u>		K-4	96.3		K-5	97.0		K-6	97.7		K-7	<u>98.1</u>	
9 5 INSPECTEUR	K-1	87.1		K-2	91.1		K-3	94.4		K-4	<u>96.0</u>		K-5	96.5		K-6	97.8		K-7	<u>98.4</u>	
0 7 DJODJO	K-1	87.0		K-2	91.0		K-3	94.2		K-4	<u>95.1</u>		K-5	95.6		K-6	96.0		K-7	96.7	
1 12 PETIT_POUCET	K-1	86.7		K-2	90.2		K-3	93.7		K-4	94.9		K-5	<u>95.3</u>		K-6	96.6		K-7	96.9	
2 4 FOOTBALL	K-1	86.4		K-2	92.6		K-3	<u>95.4</u>		K-4	97.1		K-5	97.9		K-6	99.2		K-7	99.3	
8 CHOUETTE_BOUQUET	K-1	86.4		K-2	92.3		K-3	94.8		K-4	<u>95.7</u>		K-5	96.6		K-6	98.2		K-7	98.8	
4 17 PAS_DE_SOLEIL	K-1	86.4		K-2	90.8		K-3	92.8		K-4	93.9		K-5	94.7		K-6	<u>96.5</u>		K-7	97.1	7
5 1 SOUVENIR CHERIR	K-1	86.2		K-2	91.1		K-3	92.9		K-4	<u>95.9</u>		K-5	96.8		K-6	97.1		K-7	97.4	
16 FREQUENTE_AGNEN	K-1	85.9		K-2	90.7		K-3	93.1		K-4	<u>95.6</u>		K-5	96.6		K-6	97.4		K-7	<u>98.1</u>	
7 11 JE_FUME	K-1	85.5		K-2	89.5		K-3	93.1		K-4	93.8		K-5	94.8		K-6	<u>97.0</u>		K-7	97.7	
14 MALADE	K-1	85.5		K-2	90.4		K-3	93.8		K-4	<u>95.0</u>		K-5	97.1		K-6	<u>98.0</u>		K-7	98.6	
13 VELO	K-1	84.7		K-2	92.1		K-3	93.7		K-4	94.1		K-5	<u>96.8</u>		K-6	97.6		K-7	97.9	
2 COWBOYS	K-1	83.0		K-2	89.6		K-3	91.1		K-4	93.0		K-5	94.3		K-6	<u>95.0</u>		K-7	96.4	
6 REX	K-1	82.6		K-2	90.6		K-3	92.6		K-4	93.7		K-5	93.9		K-6	94.6		K-7	<u>95.3</u>	
						1													1	1	1

This is where some facts abut French lexical distributions is useful

- Fact: Cobb and Horst (2004)
 - But why? Because of the greater number of non-optional function words?
 - Example
 - "Nicolas brought <u>a</u> ball" (1 function word, α)
 - "Nicolas <u>a</u> aporté <u>un</u> ballon" (2 function words, a and un)
- So maybe try some other Excel sorts that group the texts differently?
 - Say, by the k-levels at which 95% coverage is achieved?

So sort by Columns I, then L, then O...

A	В	С	D	E	F	G	Н	1	J	K	L	М	N	0	Р	Q	R	S	Т	U	V
1 LE PETIT NICOLAS																					
2 AS PUBLISHED																					
3 1 SOUVENIR CHERIR	K-1	86.2		K-2	91.1		K-3	92.9		K-4	<u>95.9</u>		K-5	96.8		K-6	97.1		K-7	97.4	
4 10 LOUISETTE	K-1	87.5		K-2	93.4		K-3	<u>95.3</u>		K-4	96.3		K-5	97.0		K-6	97.7		K-7	<u>98.1</u>	
5 11 JE_FUME	K-1	85.5		K-2	89.5		K-3	93.1		K-4	93.8		K-5	94.8		K-6	<u>97.0</u>		K-7	97.7	
6 12 PETIT_POUCET	K-1	86.7		K-2	90.2		K-3	93.7		K-4	94.9		K-5	<u>95.3</u>		K-6	96.6		K-7	96.9	
7 13 VELO	K-1	84.7		K-2	92.1		K-3	93.7		K-4	94.1		K-5	<u>96.8</u>		K-6	97.6		K-7	97.9	
8 14 MALADE	K-1	85.5		K-2	90.4		K-3	93.8		K-4	<u>95.0</u>		K-5	97.1		K-6	98.0		K-7	98.6	
9 15 BIEN_RIGOLE	K-1	88.2		K-2	92.3		K-3	94.4		K-4	94.9		K-5	<u>95.9</u>		K-6	96.4		K-7	96.9	
10 16 FREQUENTE_AGNEN	K-1	85.9		K-2	90.7		K-3	93.1		K-4	<u>95.6</u>		K-5	96.6		K-6	97.4		K-7	<u>98.1</u>	
11 17 PAS_DE_SOLEIL	K-1	86.4		K-2	90.8		K-3	92.8		K-4	93.9		K-5	94.7		K-6	<u>96.5</u>		K-7	97.1	
12 18 QUITTE_MAISON	K-1	88.2		K-2	92.8		K-3	94.8		K-4	<u>96.3</u>		K-5	97.1		K-6	97.9		K-7	<u>98.4</u>	
13 2 COWBOYS	K-1	83.0		K-2	89.6		K-3	91.1		K-4	93.0		K-5	94.3		K-6	<u>95.0</u>		K-7	96.4	
14 3 BOUILLON	K-1	88.2		K-2	92.4		K-3	94.2		K-4	94.9		K-5	<u>95.7</u>		K-6	96.5		K-7	96.8	
15 4 FOOTBALL	K-1	86.4		K-2	92.6		K-3	<u>95.4</u>		K-4	97.1		K-5	97.9		K-6	99.2		K-7	99.3	
16 5 INSPECTEUR	K-1	87.1		K-2	91.1		K-3	94.4		K-4	<u>96.0</u>		K-5	96.5		K-6	97.8		K-7	<u>98.4</u>	
17 6 REX	K-1	82.6		K-2	90.6		K-3	92.6		K-4	93.7		K-5	93.9		K-6	94.6		K-7	<u>95.3</u>	
18 7 DJODJO	K-1	87.0		K-2	91.0		K-3	94.2		K-4	<u>95.1</u>		K-5	95.6		K-6	96.0		K-7	96.7	
19 8 CHOUETTE_BOUQUET	K-1	86.4		K-2	92.3		K-3	94.8		K-4	<u>95.7</u>		K-5	96.6		K-6	98.2		K-7	98.8	
20 9 CARNET	K-1	88.6		K-2	93.8		K-3	<u>95.1</u>		K-4	95.5		K-5	97.6		K-6	98.4		K-7	98.7	
21																					

Voilà – a good first guess at a set of four distinct lexical levels across these 18 stories

	•		Ju		V C		4 ()	CI	. 00	<u> </u>			 . •					
4:	3																	
4																		
	SORTED BY	%	16.4	00.4		14.0	00.0	16.0	05.4		16.4	07.4	14.5	07.0	16.0	00.0	14.7	00.0
4	4 FOOTBALL		K-1	86.4		K-2	92.6	K-3	<u>95.4</u>		K-4	97.1	K-5	97.9	K-6	<u>99.2</u>	K-7	99.3
4	7 10 LOUISETT		K-1	87.5		K-2	93.4	K-3	<u>95.3</u>		K-4	96.3	K-5	97.0	K-6	97.7	K-7	<u>98.1</u>
48	9 CARNET		K-1	88.6		K-2	93.8	K-3	<u>95.1</u>		K-/	95.5	K-5	97.6	K-6	<u>98.4</u>	K-7	98.7
49	18 QUITTE_M	ISON	K-1	88.2		K-2	92.8	K-3	94.8		λ-4	<u>96.3</u>	K-5	97.1	K-6	97.9	K-7	<u>98.4</u>
50	5 INSPECTEUR		K-1	87.1		K-2	91.1	N. 9	94		K-4	<u>96.0</u>	K-5	96.5	K-6	97.8	K-7	<u>98.4</u>
5	1 SOUVENIR C	IERIR	K-1	86.2		K-2	91.1	K-3	92.9		K-4	<u>95.9</u>	K-5	96.8	K-6	97.1	K-7	97.4
5	2 8 CHOUETTE_	DUQUET	K-1	86.4		K-2	92.3	K-3	94.8		K-4	<u>95.7</u>	K-5	96.6	K-6	<u>98.2</u>	K-7	98.8
5	16 FREQUENT	_AGNEN	K-1	85.9		K-2	90.7	K-3	93.1		K-4	<u>95.6</u>	K-5	96.6	K-6	97.4	K-7	<u>98.1</u>
54	4 7 DJODJO		K-1	87.0		K-2	91.0	K-3	94.2		K-4	<u>95.1</u>	K-5	90.0	K-6	96.0	K-7	96.7
5	14 MALADE		K-1	85.5		K-2	90.4	K-3	93.8		K-4	<u>95.0</u>	y 5	97.1	K-6	<u>98.0</u>	K-7	98.6
56	5 13 VELO		K-1	84.7		K-2	92.1	K-3	93.7		K 1	94.1	K-5	<u>96.8</u>	K-6	97.6	K-7	97.9
5	7 15 BIEN_RIG	LE	K-1	88.2		K-2	92.3	K-3	94.4		K-4	94.9	K-5	<u>95.9</u>	K-6	96.4	K-7	96.9
58	3 BOUILLON		K-1	88.2		K-2	92.4	K-3	94.2		K-4	94.9	K-5	<u>95.7</u>	K-6	98	K-7	96.8
59	12 PETIT_PC	CET	K-1	86.7		K-2	90.2	K-3	93.7		K-4	94.9	7-5	<u>95.3</u>	K-6	96.6	K-7	96.9
60	11 JE_FUM		K-1	85.5		K-2	89.5	K-3	93.1		K-4	93.8	K	94.8	K-6	<u>97.0</u>	K-7	97.7
5	1 17 PAS_D/_S0	DLEIL	K-1	86.4		K-2	90.8	K-3	92.8		K-4	93.9	K-5	94.7	K-6	<u>96.5</u>	K-7	97.1
	2 COWB(/S		K-1	83.0		K-2	89.6	K-3	91.1		K-4	93.0	K-5	94.3	K-6	<u>95.0</u>	K-7	96.4
6	6 REX		K-1	82.6		K-2	90.6	K-3	92.6		K-4	93.7	K-5	93.9	K-6	94.6	K-7	<u>95.3</u>
6	4																	

This would need empirical evaluation Ongoing study design:

- Group A reads the stories in published sequence
- Group B reads the stories in the DDL-revised sequence
 - → Are there differences in rate, comprehension, enjoyment, lexical development, grammar development?
 - → Is there impact on writing (sentence length, accuracy, VP?)

Perhaps the point is made

DDL has a lot to offer in the design of this course Beyond just 'choosing a good book the learners will like'

- But there are several further piecemeal ways DDL can operate in this course
 - Corpus searches of the stories by teacher to explore and teach features of the stories
 - Is a particular word recurring or one-off?
 - Control of examination lexis
 - (Do questions and rubrics about the stories match the language of the stories?)
 - If a corpus of teacher-talk is available, are teachers reinforcing the language of the texts in the language of the classroom?
 - Is there reinforcement and recycling of new words?

Par contre ~

13 le velo=0 16 frequent agnan=0

RANGE: apercevoir appears in 2 out of 14 sub-corpora (14.29%)

```
Home > Concordancers > French Input > Output
                                                   (« Back keeps original settings)
                                                                                Copiable extract-link to this data >> here
Concordance for (f) lemma apercevoir in petit nicolas Dictionnaire/TTS Free Eng v
   Extract checked ↓ items: ✓ All | Ø | any10 | 20 | 30 | 50 Go >
  RANGE: apercevoir appears in 2 out of 14 sub-corpora (14.29%)
     Click any KEYWORD for more context + 'Get More' option FEB 2022
   3 hits
  001. otre vie, c'est raté, parce qu'on s'est APERÇU que le photographe n'était plus l 01_souve...
  002. quoi? a demandé Alceste. Je viens de m'APERCEVOIR que nous n'avons pas de ballo 04_le_fo...
  003. amais, comme aujourd'hui, je ne me suis APERÇU à quel point notre métier est un 05_1'ins...
    Distribution for lemma apercevoir: 01 souvenir cherir=1 04 le football=1 05 linspecteur=1 02 cow-bo
    05 l'inspecteur=0 06 rex=0 07 djodjo=0 08 chouette bouquet=0 09 les carnets=0 10 louisette=0 11 je fu
```

T can go as deep as the Ss' needs, e.g. collocation

```
Home > Concordancers > French Input > Output (« Back keeps original settings) Copiable extract-link to this data >> h

Concordance for (f)lemma répondre in petit_nicolas Dictionnaire/TTS Fren_Eng ▼

Extract checked ↓ items: ✓ All | Ø | any10 | 20 | 30 | 50 Go >
```

RANGE: répondre with que on right side appears in 8 out of 14 sub-corpora (57.14%)

Click any KEYWORD for more context + 'Get More' option FEB 2022

```
10 hits

001.  plaisir à votre maîtresse? » Nous avons RÉPONDU que oui, parce que nous l'aimons 01_souve...

002.  dit de cesser de manger, mais Alceste a RÉPONDU qu'il fallait bien qu'il se nour 01_souve...

003.  pachim? » a demandé Clotaire. Joachim a RÉPONDU qu'il ne voyait rien. Alors, Clo 03_le_bo...

004.  ce qu'il lui était arrivé et il nous a RÉPONDU qu'il s'était endormi à force de 03_le_bo...

005.  mauvais côté du terrain! Moi, je lui ai RÉPONDU que si le soleil ne lui plaisait 04_le_fo...

006.  uérir, mais ça fait très mal. Moi, j'ai RÉPONDU à papa que Rex n'était pas malad 06_rex

007.  ndé, Eudes. « Je ne peux pas, je lui ai RÉPONDU, il faut que je rentre chez moi 08_choue...

008.  re semblant de chanter et Alceste lui a RÉPONDU que il ne faisait pas semblant de 10_louis...

009.  as du feu? » il m'a demandé, je lui ai RÉPONDU que non. « Ben alors, a dit Alce 11_je_fu...

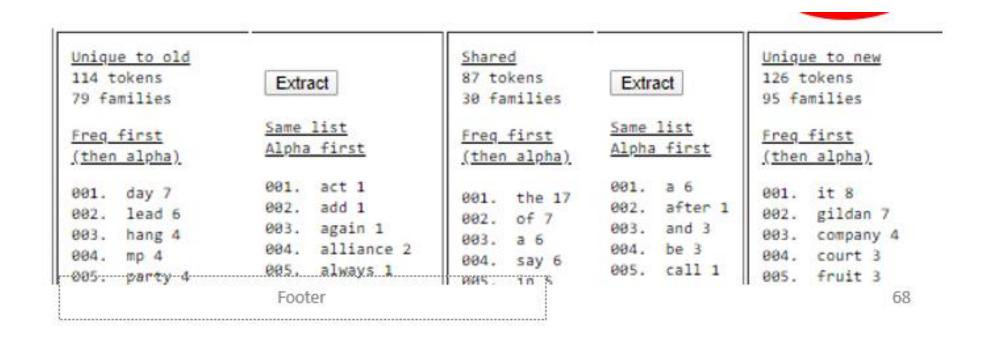
010.  n me ferait mettre en prison. Je lui ai RÉPONDU que j'avais bien envie de lui ta 16_frequ...
```

E.g., syntax



Plus two uses for Lextutor's Text_Lex_Compare Text-lexis comparison

Program determines the words/lemmas shared between two texts vs.
 unique to one text



1. If you can get a corpus of teacher talk for this course...

- Use Text_Lex_Compare to check how much lexis
 of texts is getting recycled in the classroom
- And if little is recycled, then include re-use/recycle practice in TT (teacher training) programs
 - Low transfer is typical →

Home > Text Lex Compare

Text Lex Compare v4

Research

TEXT COM

· Subtract one text from another; one list from another; a list from a text

- Find the degree of word repetition of chapters, etc)
- Researchers Find the coverage of
- Teachers See what's new and wh

Method A] Old or first text here

Alliance MP calls on Day to resign desp Canadian Alliance Leader Stockwell Day dissent over his leadership Wednesday a unite his caucus, with one MP again cal Despite Day's claim that he had unanim FAMILY ANALYSIS

Titles: Old COURSE UNIT

Empty Count

Home >> INPUTS Via Back Direct >> Text Lex Output

New words in second text

I Index-Edit-Area at bottom

First text: COURSE UNIT (109 families)

Second text(s): TEACHER TALK (125 families)



plan, Art Hanger emerged from a private Using the family as unit of comparison means that if cat is in Text 1 and cats in Text 2 then this is considered a repetit Note also that in this routine, "unfamilied" items revert to classification by type (e.g., repeated proper nouns)

> TOKEN Recycling Index: (87 repeated tokens : 213 tokens in new text) = 40.85% FAMILIES Recycling Index: (30 repeated families: 125 families in new ext) = 24.00%

(Token recycling will normally be the most interesting measure of e.g. coverage and hence text completensibility, as with y

<u>Unique to old</u> 114 tokens 79 families	Extract
Freq first	<u>Same list</u>
(then alpha)	<u>Alpha first</u>
001. day 7	001. act 1
002. lead 6	002. add 1
003. hang 4	003. again 1
004. mp 4	004. alliance 2
005. party 4	005. always 1

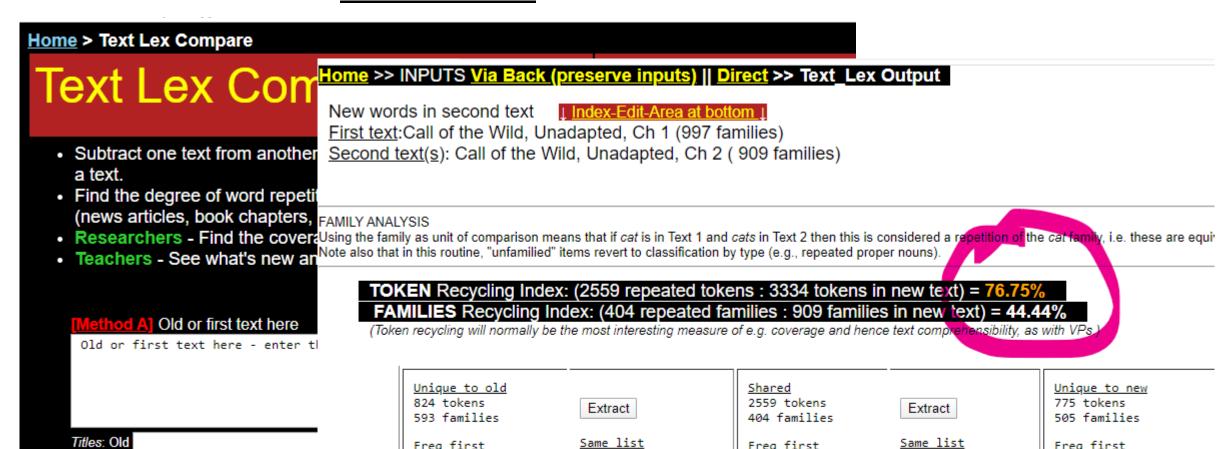
<u>Shared</u> 87 tokens 30 families	Extract
Freq first (then alpha)	<u>Same list</u> <u>Alpha first</u>
001. the 17 002. of 7 003. a 6 004. say 6	001. a 6 002. after 1 003. and 3 004. be 3 005. call 1

Unique to new 126 tokens 95 families Freq first (then alpha) 991. it 8 002. gildan 7 003. company 4 004. court 3 005. fruit 3

2. And finally EXAMS

To write fair examination texts/questions

 Use Text_Lex_Compare to assure that 95% of the words on the test were even seen in the course



+ Let VP identify words to gloss in an exam

ilers > VP-Compleat Input > Output (Use 'Back' to preserve previous inputs) FRAM

colleague_[1] communicate_[1] complaint_[2] compromise_[1] conclude_[1] concongress_[2] consistent_[1] constitute_[1] coordinate_[2] counsel_[2] criticise_[1] defend_[8] delay_[3] despite_[1] diplomat_[2] disrupt_[2] echo_[1] effective_[1] e highlight_[1] initial_[2] inquire_[11] intelligence_[4] interfere_[1] internal_[1] internal_[1] poll_[1] portion_[1] praise_[1] probe_[1] procequest_[1] respond_[2] response_[2] resume_[1] reverse_[1] scandal_[2] senate substance_[1] survey_[1] target_[1] task_[1] veteran_[1]

BNC-COCA-K4 Families: [fams 11 : types 11 : tokens 22]

VP-negative: bnc_coca-4

ambassador_[2]_shaos_[1] coherent_[1] credible_[1] opt_[1] rocked_[1] senator_testimony_[10] undermine_[1]

BNC-COCA-K5 Families: [fams 14: types 16: tokens 20]

VP-negative: bnc coca-5

aide_[4] amplify_[1] butt_[1] convene_[1] delete_[1] electorate_[1] escalate_[2] sway_[1] testify_[3] vicious_[1] withhold_[1]

EXAM QUESTION

Text text text. Text text text
Text text text. Text text text
undermine* text text text
text...

* Undermine – try to secretly weaken someone's plan or strategy

Or 2% left to inference <u>if</u> that was part of the course

And teacher training?

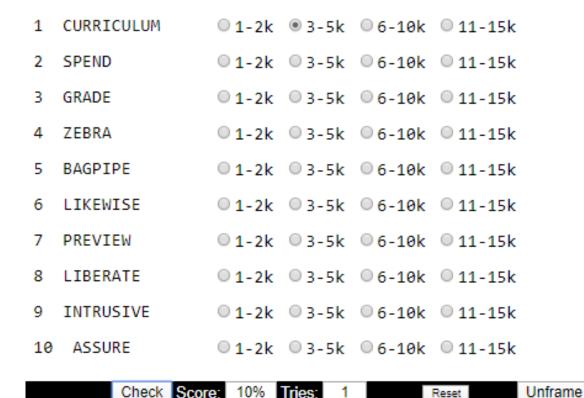
Freq Train V.2

Frequency intuitions trainer - BNC or BNC-COCA frequency lists

OK

- Teach use of DDCD tools?
- T-Training AT LEAST should include more needs analysis
- But, even more, AWARENESS RAISING in word frequency
 - Learners feel frequency effects,
 NSs do not

Test your Frequency Intuitions: Guess the band for each word Frequency Scheme : bnc_lists



So we have seen some pieces of researchsupported DDL – that can be put together usefully

LEXTUTOR ROUTINES INVOLVED

- TESTS
- VP / VP-Cloze
- RANGE
- TL_COMPARE
- GROUP LEX All @ www.lextutor.ca/ xxx

Many of these ideas are validated

- Boulton & Cobb (2017)
 - DDL for learning
- Cobb Horst & Nicolae (2005)
 - Group Lex
- Laufer & Ravenhorst (2010)
 - the 95% and 98% cut-offs
- Cobb (2020) 'Corpus for courses'
 - A published earlier version of this presentation

QUESTIONS? COMMENTS? MORE LIVE DEMOS?

Further reading on this topic

Earlier version:
2021 conference paper
at Swiss Conference of
Applied Linguistics

Download at **Lextutor.ca/cv/**

Corpus for courses: Data-driven course design

Thomas COBB

Université du Québec à Montréal (UQAM)
Faculté des sciences de l'éducation
1205, rue Saint-Denis, Montréal (Québec) H2X 3R9, Canada cobb.tom@ugam.ca; www.lextutor.ca

L'apprentissage sur corpus, en tant qu'ensemble de principes et de technologies, joue un rôle bien établi dans l'apprentissage des langues, et cet article montre comment ceux-ci peuvent s'appliquer également à la conception de cours ou de curriculum de langue. Si les supports de cours sont reformulés sous forme de corpus, on a alors à disposition un certain nombre de moyens pour amener la recherche sur l'apprentissage directement dans la salle de classe. Cet article commence par une définition des termes, un examen de la place de l'apprentissage sur corpus dans l'acquisition de la langue, et montre des moyens d'appliquer cette approche à la conception, au testing, et à l'évaluation de l'enseignement des langues. Le contexte est la refonte d'un cours de lecture en cours pour adultes en langue seconde basé sur une collection de documents authentiques trouvés sur Internet. Les principes et technologies de la conception de cours sur corpus ont été utilisés, premièrement, pour exposer les faiblesses de cette collection de matériels en tant que cours, et, deuxièmement, pour montrer des moyens concrets de les améliorer. Tous les logiciels impliqués dans ce travail sont accessibles au public.

Mots-clés:

apprentissage basé sur les consultations d'un corpus, conception pédagogique, analyse de besoins, vocabulaire.

Keywords:

corpus-based learning, data-driven learning, instructional design, needs analysis, text analysis, vocabulary.

1. Background & proposition

Data-driven learning (DDL) is an input and comprehension based approach to