

Language Learning: A Data-Driven Approach



Michael C. Frank
Stanford University

Intros

- Michael Frank (call me Mike)
 - Faculty at Stanford
 - Study language learning, reproducibility and replicability, pragmatics, development, metascience
 - Fun fact: play the mandolin in a bluegrass band
- You:
 - Institution, program of study
 - Research interest
 - Fun fact – or boring fact if you feel stressed

Learning goals

- Discuss the "standard model" framework for early word learning, focusing on input, processing, and uptake constructs,
- Compare different instruments and approaches for measuring child language,
- Learn a reproducible workflow for exploring language acquisition data in R, and
- Explore data from Wordbank, CHILDES, and Peekbank as a source of insights into language learning.

Course outline

- Today: Introduction to course framework and toolset
- Tuesday: Digging deeper into vocabulary data in
Wordbank
- Wednesday: exploring transcripts of children's speech
through **CHILDES** and **childe-db**
- Thursday: analyzing eye-tracking data using
Peekbank
- Friday: Student-led mini-projects

Course prerequisites

- Minimal **github** to synchronize most recent course files
- **R** and **R studio** to manipulate data
- Tidyverse including **dplyr** and **ggplot** to manipulate and visualize data
- All of these will be (briefly) introduced
- But if you do not have some of these, the second part of most classes will be frustrating

Readings and other homework

- There will be two articles assigned for each class
 - Day 1
 - Frank et al. (2021) Chapter 1, "Theoretical Foundations"
 - Kachergis, Marchman, & Frank (2022), "Towards a standard model"
 - Day 2
 - Bates & Goodman (1997), "On the inseparability of grammar and the lexicon"
 - Frank et al. (2021), Chapter 13, "Morphology, Grammar, and the Lexicon"
 - Day 3
 - MacWhinney & Snow (1990), "CHILDES: an update"
 - Sanchez*, Meylan* et al. (2019), "childe-db: a flexible and reproducible interface"
 - Day 4
 - Fernald et al. (1998), "Rapid gains in speed of verbal processing"
 - Zettersten et al. (2022), "Peekbank: an open, large-scale repository"
- Beyond these, please work to complete the Rmds from class
- On Day 4, you should send me an email with your group and your proposed project (1 paragraph)

Outline

- 1. Basics of early word learning**
2. Toward a “standard model” of language learning
3. Developing resources (to measure quantities in the model)
 - Input – childe-db
 - Learning – wordbank
 - Processing – peekbank
4. The course toolset



Mary Cassatt - The Complete Works



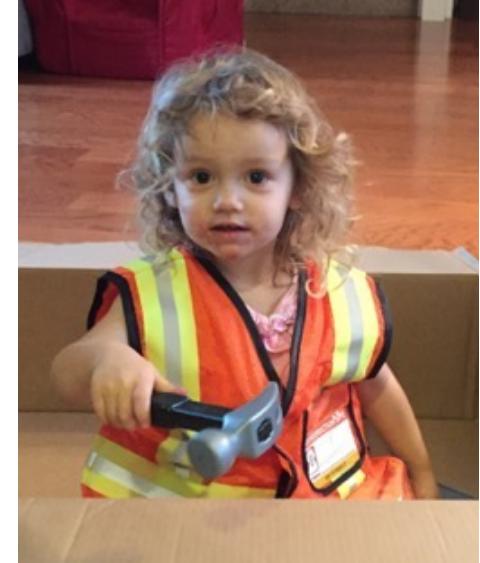
An explosion of language

18 mo: “happy-b”

19 mo: “blue ball”

23 mo: “spike doggy no food
eat dirt”

26 mo: “dada move own body,
my need lilbit more
space”



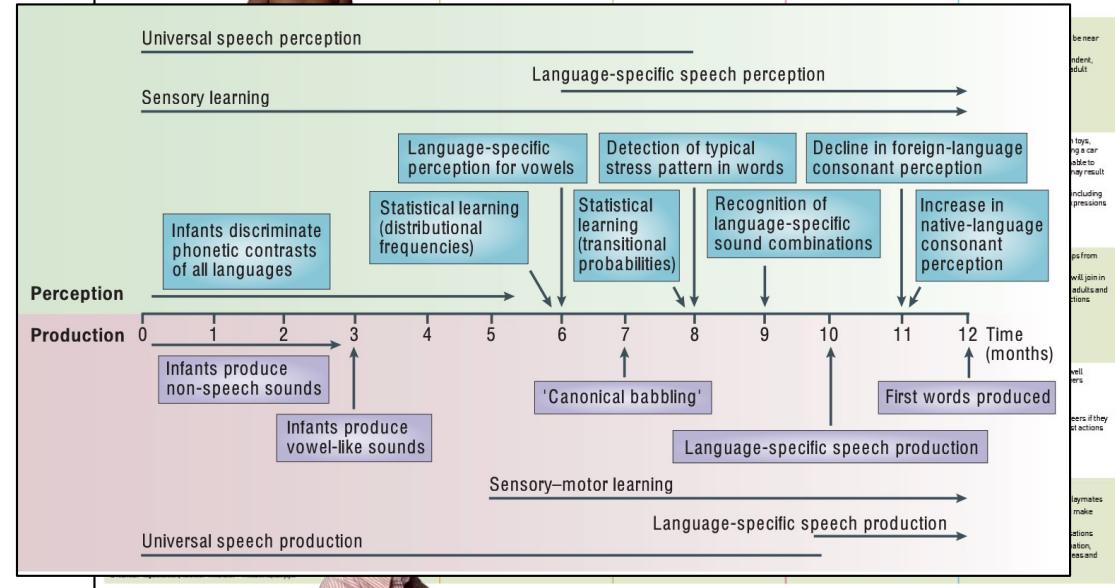
Ages and stages theories are not sufficient

A guide for early years practitioners



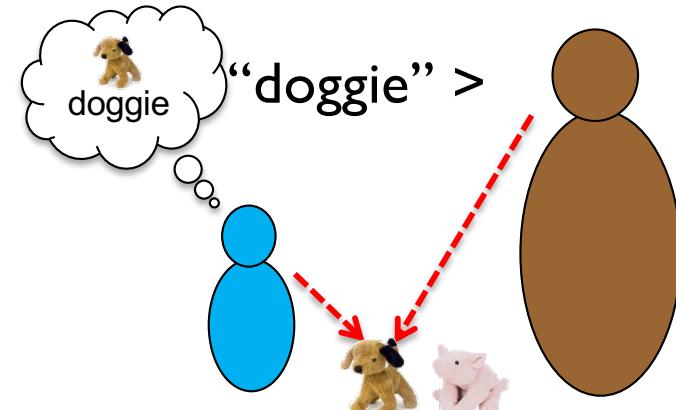
Stages of Speech and Language Development

	Listening and attention	Understanding	Speech sounds and talk	Social skills
Up to 3 months	<ul style="list-style-type: none"> Turns towards a familiar sound Started by loud noises 	<ul style="list-style-type: none"> Recognises parent's voice Often calmed by familiar friendly voice, e.g. parent's 	<ul style="list-style-type: none"> Frequently cries especially when unhappy or uncomfortable Makes vocal sounds, e.g. cooing, gurgling 	<ul style="list-style-type: none"> Looks at faces and copies facial movements, e.g. sticking out tongue! Makes eye contact for fairly long periods
3 – 6 months	<ul style="list-style-type: none"> Watches face when someone talks 	<ul style="list-style-type: none"> Show excitement at sound of approaching voices 	<ul style="list-style-type: none"> Makes vocal noises to get attention Makes sounds back when talked to Laughs during play Babbles to self 	<ul style="list-style-type: none"> Shows different emotions in parent's voice and may respond differently, for example, smile, quieten, laugh Cries in different ways to express different needs
6 – 12 months	<ul style="list-style-type: none"> Locates source of voice with accuracy Focuses on different sounds, e.g. telephone, doorbell, clock 	<ul style="list-style-type: none"> Understands frequently used words, such as: 'all gone', 'milk' and 'bye-bye' Stops and looks when hears own name Understands simple instructions when supported by gestures and context 	<ul style="list-style-type: none"> Uses speech sounds (babbling) to communicate with adults, says sounds like 'ba-ba, no-no, ga-ga' Stops babbling when hears familiar adult's voice Uses gestures such as waving and pointing to help communicate Around 12 months begins to use single words e.g. 'mum-mum', 'sad', 'see teddy' 	<ul style="list-style-type: none"> Enjoys action rhymes and songs Reacts to cosy adult speech and lip movements Takes 'turns' in conversations (using babble)
12 – 15 months	<ul style="list-style-type: none"> Attends to music and singing Enjoys sound-making toys/objects 	<ul style="list-style-type: none"> Understands single words in context, e.g. car, milk, daddy Understands more words than they can say Understands single instructions, e.g. 'kiss mummy', 'give to daddy', 'stop' 	<ul style="list-style-type: none"> Say around 10 single words, although these may not be clear Reaches or points to something they want while making speech sounds 	<ul style="list-style-type: none"> Likes being with familiar adults Likes watching adults for short periods of time



Early vocabulary as a case study

- Easy to observe and measure
- Units (words) are less controversial than representations underlying morphosyntax
- Nonetheless tightly linked to morphosyntax (day 2)
- Grounded in concept learning, social cognition, etc.



Language input

	Yearly total				Yearly CDS			
	Hours		Words (M)		Hours		Words (M)	
Urban, high SES								
H&R (N=13) ^t	1221 ^{w,c}	[578, 1987]	11.0 ^c	[5.20, 17.9]	786 ^w	[372, 1279]	7.07	[3.35, 11.5]
S&G (N=6) ^t	2023 ^{w,m}	[1243, 2858]	18.2 ^m	[11.2, 25.7]	1223 ^{w,m}	[853, 1574]	11.0 ^m	[7.7, 14.2]
VdW (N=1)	931		9.28		140		1.39	
Urban, low SES								
H&R (N=6) ^t	363 ^{w,d}	[136, 558]	3.26 ^d	[1.22, 5.02]	225 ^w	[84, 346]	2.02	[0.76., 3.11]
W&F (N=29) ^t	363 ^w	[52, 1049]	3.27	[0.46., 9.44]	225 ^w	[32, 650]	2.03	[0.29, 5.85]
Rural, low SES								
S&G (N=6) ^t	503 ^{w,m}	[365, 640]	4.53 ^m	[3.28, 5.76]	234 ^{w,m}	[132, 322]	2.10 ^m	[1.19, 2.90]

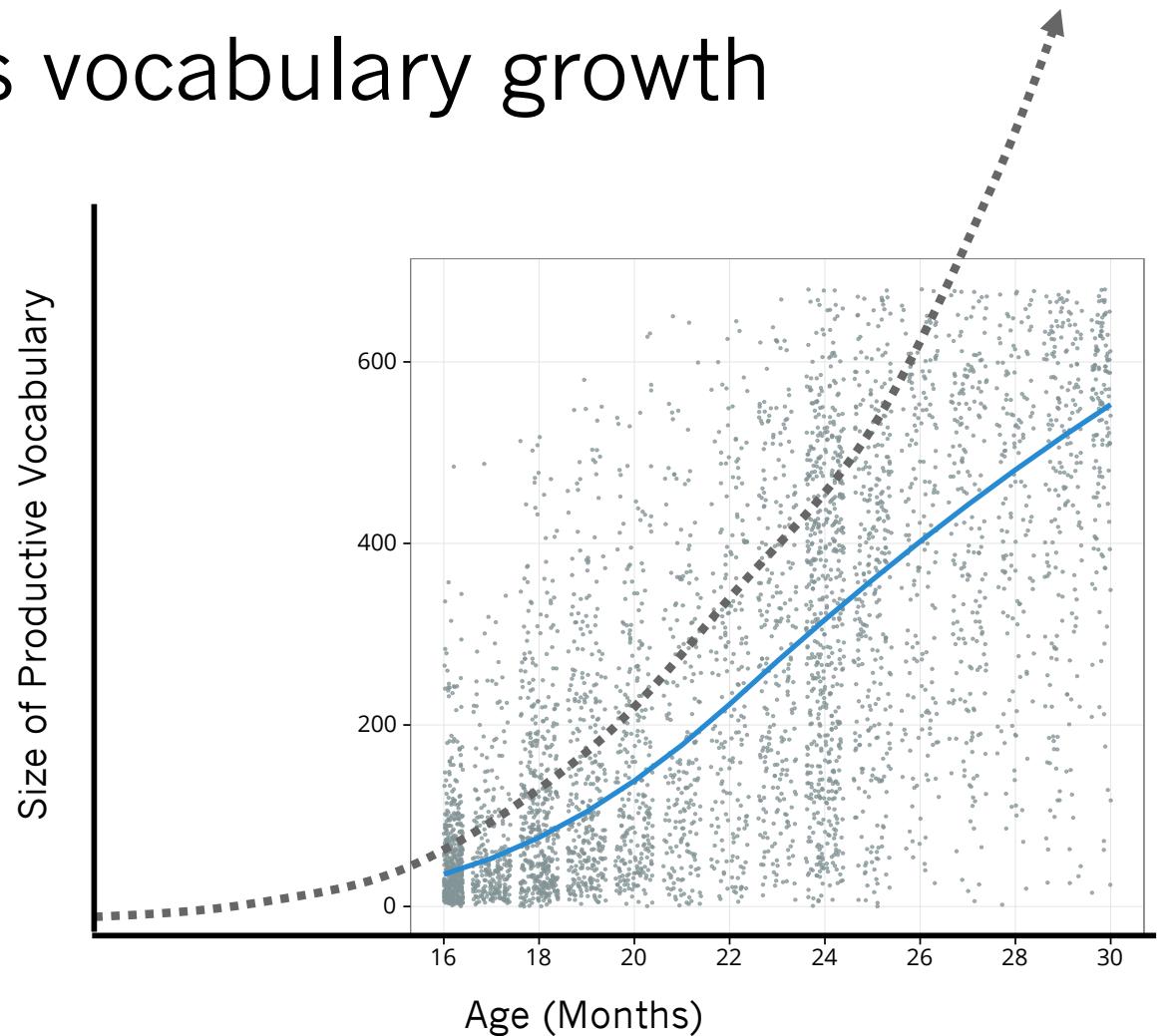
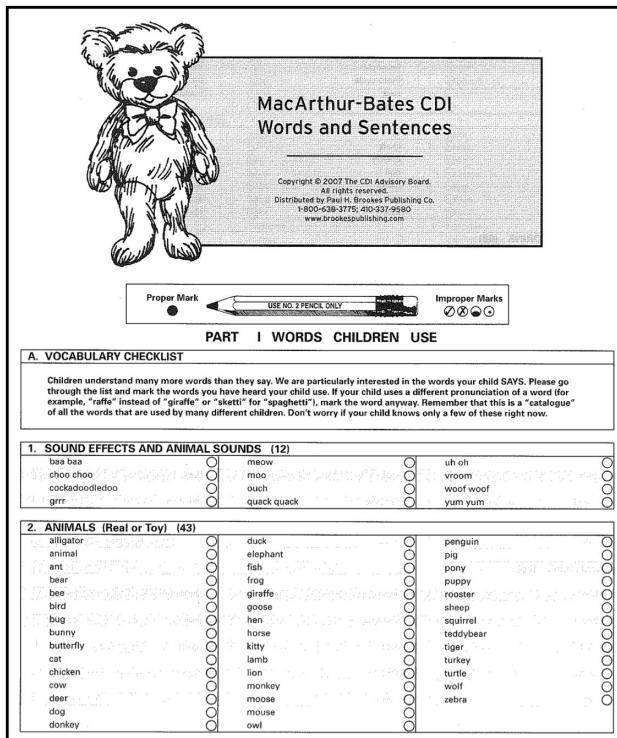
Surprisingly tricky to get accurate counts of words per year (need transcripts...)

All speech: 250k - 1.5M words/month

Child-directed: 150k – 1M words/month

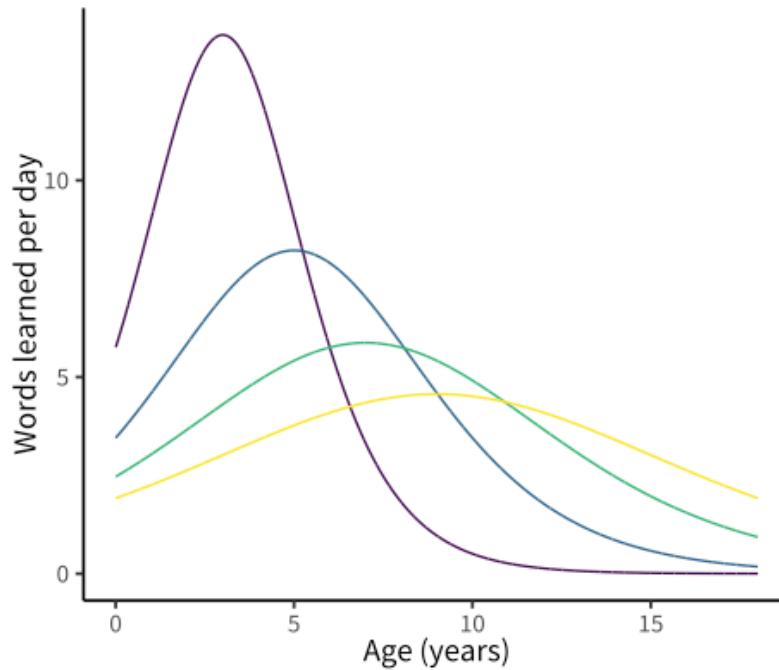
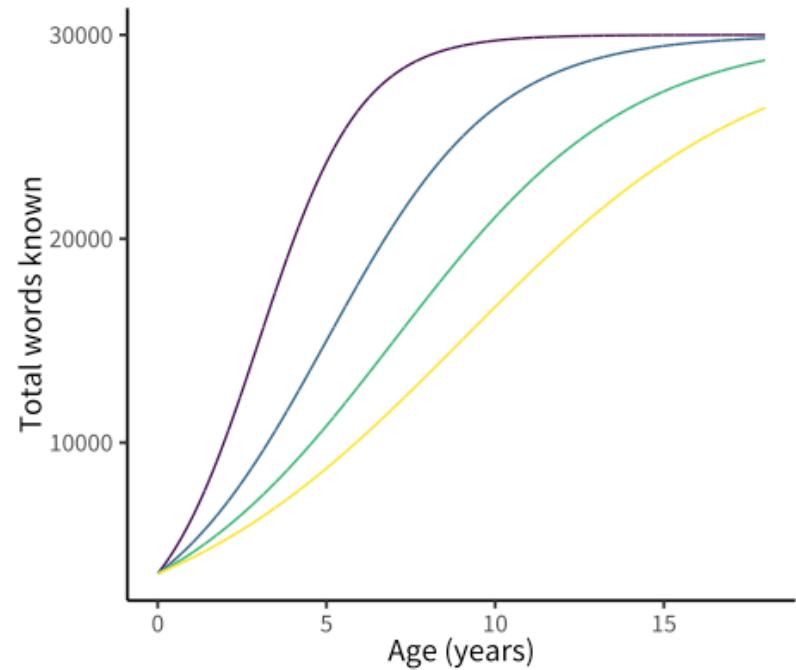
Dupoux (2018)

Children's vocabulary growth



CDI data at <http://wordbank.stanford.edu>

The rate of word learning



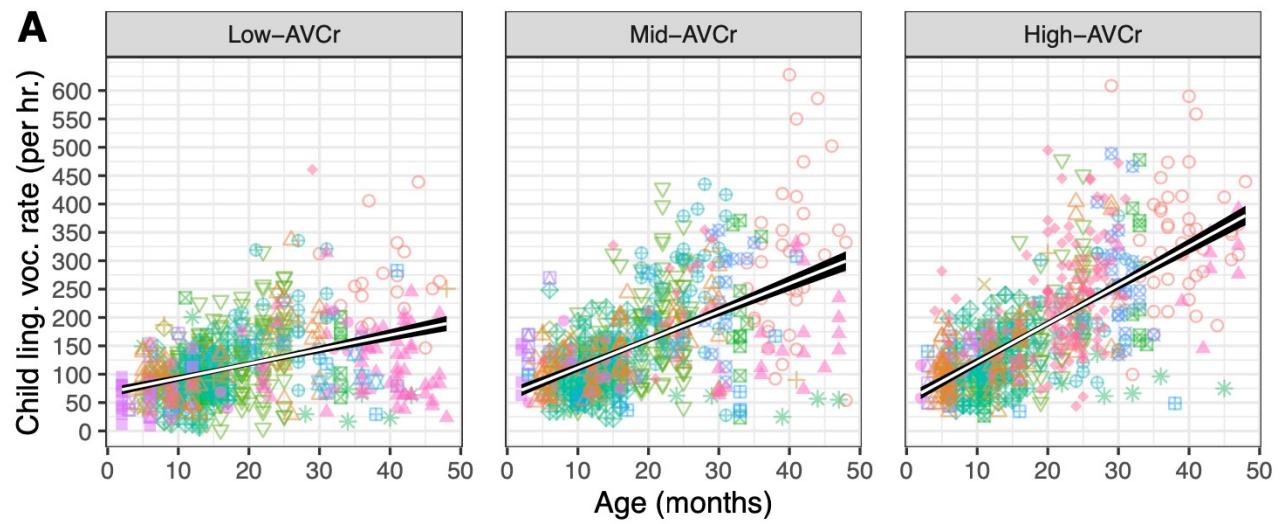
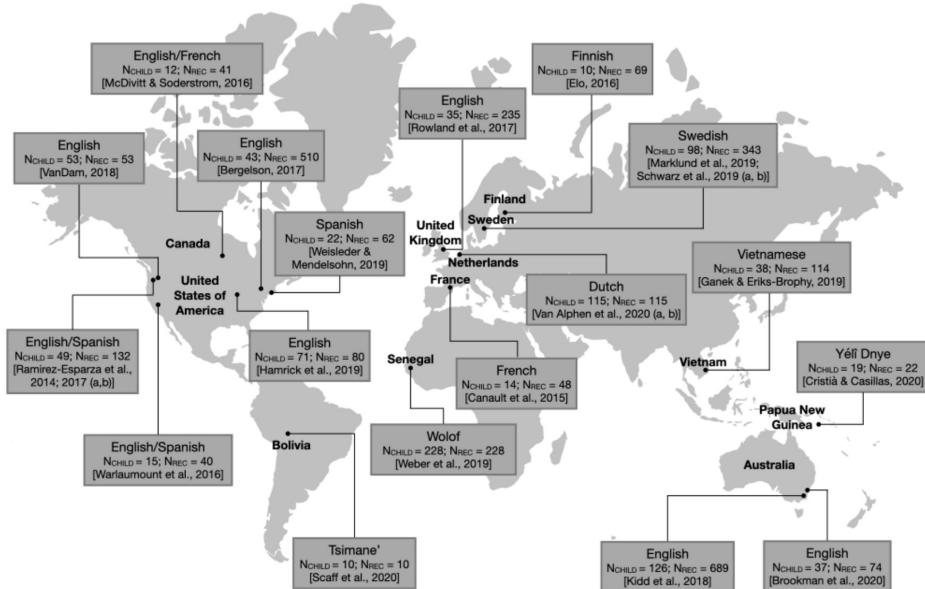
Children are “learning” several words each day!

Input to uptake

Automated analysis of vocalizations through 1000 daylong audio recordings reveals robust input-uptake correlations



Bergelson et al. (psyarxiv)



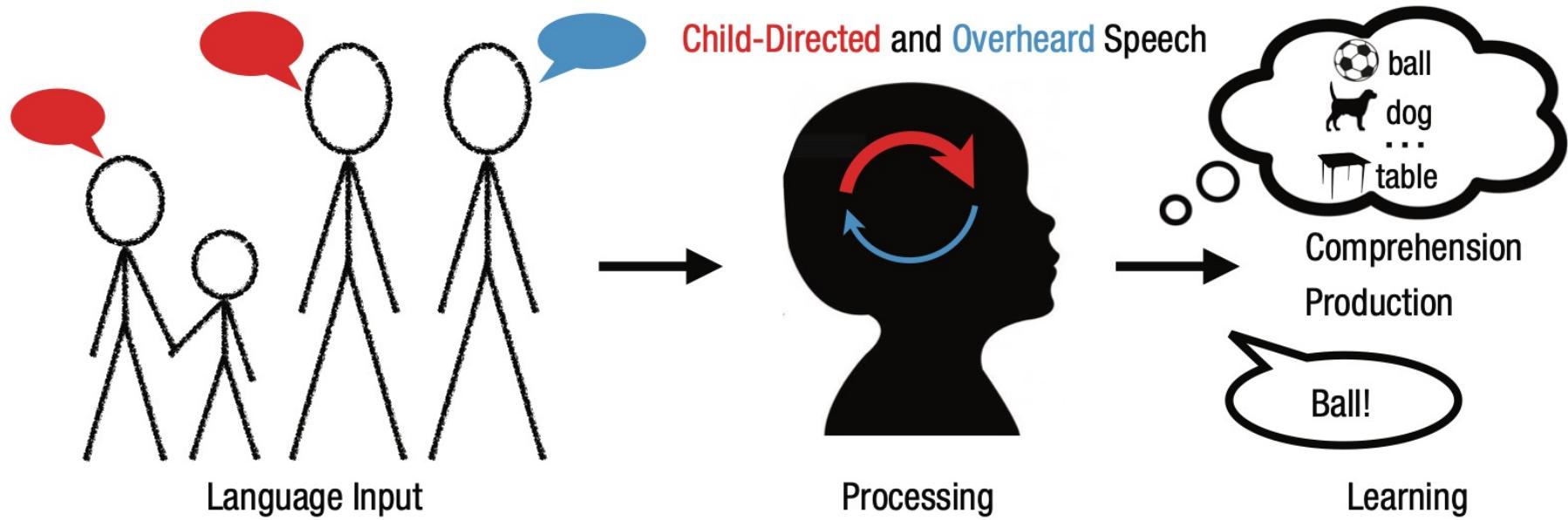
Outline

1. Basics of early word learning
2. **Toward a “standard model” of language learning**
3. Developing resources (to measure quantities in the model)
 - Input – childe-s-db
 - Learning – wordbank
 - Processing – peekbank
4. The course toolset

A “standard model” for language learning



Everyone has (more or less) the same theory in mind



Can we formalize this theory and put it in contact with data? What data?

Kachergis, Marchman, & Frank (2022), *Current Directions in Psych Sci*

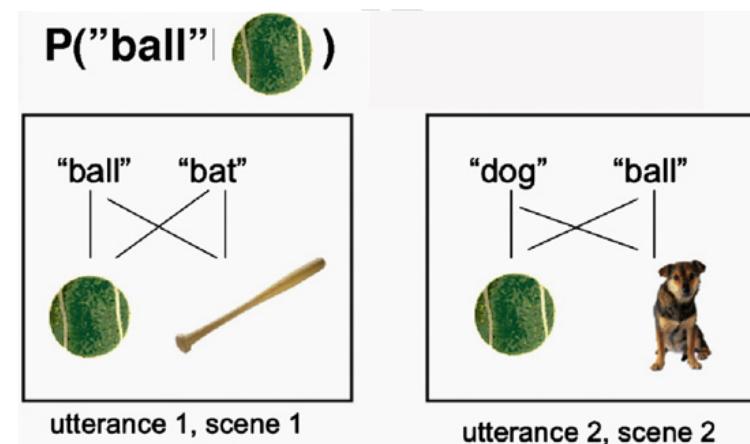
Associative learning foundations

"Statistical learning"

golabupadotitupiro
bidakugolabupadoti
golabupadotitupiro
bidakugolabupadoti

Saffran, Aslin, & Newport (1996)

"Cross-situational learning"



Yu & Smith (2008)

Defusing the Childhood Vocabulary Explosion

Bob McMurray

Between birth and adulthood, children learn about 60,000 words, on average, 8 to 10 words per day. Studies consistently reveal that, during the second postnatal year, word learning accelerates dramatically (1). Although the acceleration is continuous and not stagelike (2, 3), this so-called vocabulary explosion is a foundational phenomenon that theories of language acquisition must address.

acquisition threshold, it is learned. This model exhibits the characteristic pattern of slow learning followed by acceleration (Fig. 1B, black line). Deceleration is seen at the end of acquisition, something that has been hypothesized (3) but not examined empirically [Supporting Online Material (SOM) text S2].

Further simulations examined mechanisms that leverage initial words to facilitate learning

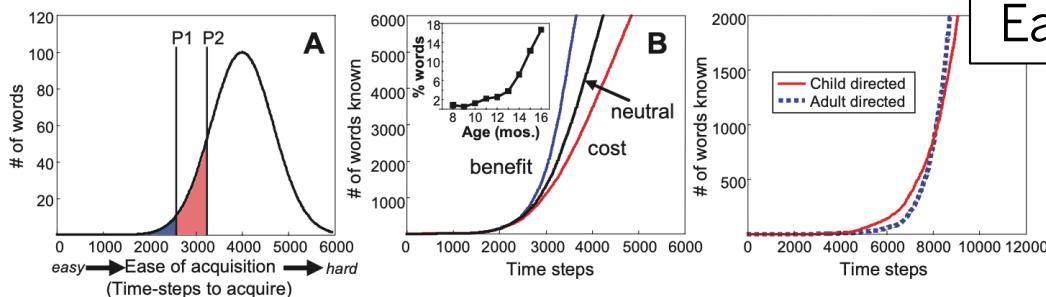


Fig. 1. (A) A Gaussian distribution of time to acquisition. Between 0 and 2600 time steps (P1), 1481 words are acquired. An additional 3966 are acquired in the next 600 (P2). (B) Vocabulary size as a function of time in initial simulations, and when learning a word offers a cost or benefit to future learning. (Inset) Percentage of words on MacArthur Communicative Development Inventory produced by children as a function of age (SOM text S1). (C) Acquisition in simulations based on word frequency.

Language Input

ball

dog

...

table

dog

have

Accumulation

Learned →



Each word is a bucket, filled by input

processes are not causally necessary to explain acceleration; such processes may, in fact, arise in response to acceleration, or they may offset other processes that slow learning (6). Moreover, the model's generality suggests that any parallel learning system (e.g., motor patterns and concepts) should behave similarly. Acceleration is an unavoidable by-product of variation in difficulty. It should not be misconstrued as evidence for specialized learning.

On Leveraged Learning in I Relationship to

Colleen Mitchell,^a

^a*Department of Mathematics, and*
^b*Department of Psychology, and*

Received 10 September 2008; received in revised

Abstract

Children at about age 18 months experience acquisition is a robust phenomenon, although the exact shape of explanations, which we term collectively as lever words helps with the learning of others. In this framework, learning is slow. As more words are acquired, new words that fuel the explosion in learning. learning in the vocabulary spurt by proposing a simple model that leverage can change both the shape and timing of the spurt. If leverage did not exist, the vocabulary spurt would be a single, sharp peak. However, if it did not exist in the corresponding model, the Zipfian distribution of word frequencies would be violated, but this is not the case. The distribution is complex.

Can these

Keywords:

Can these models be used to connect variability in input to variability in uptake?

OPEN  ACCESS Freely available online



A Computational Model Associating Learning Process, Word Attributes, and Age of Acquisition

Shohei Hidaka*

Japan Advanced Institute of Science and Technology

Abstract

We propose a new model-based computational tool for understanding AoAs as measures of vocabulary growth between three theoretical factors: how different learning processes, required to acquire a word, likely affect it in three respects. The first analysis shows that the standard alternative. The second analysis shows that psychological attributes, such as frequency, are a significant trend predicted by our estimated model on word learning in children. We

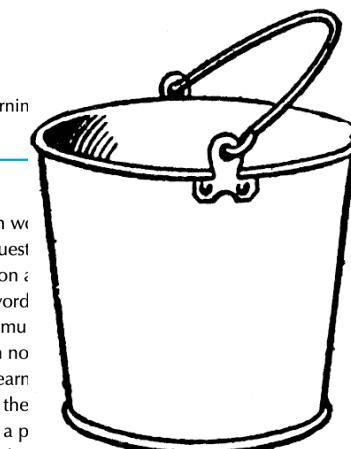
Citation: Hidaka S (2013) A Computational M
journal.pone.0076242

Editor: Johan J. Bolhuis, Utrecht University, The Netherlands



How Data Drive Early Word Learning: A Cross-Linguistic Waiting Time Analysis

Francis Mollica¹ and Steven T. Piantadosi¹



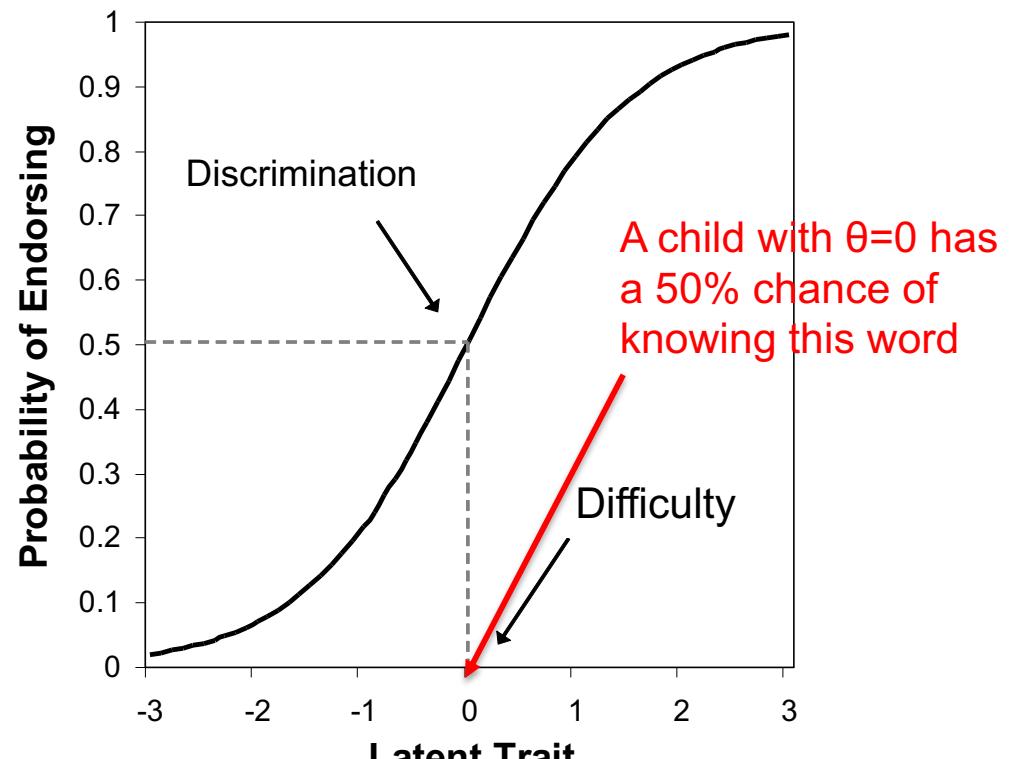
Keywords: word learnin

ABSTRACT

The extent to which words influence learning on a word-by-word basis can rapidly learn word information across multiple word acquisition nodes in early language learning. The model accurately predicts the parameters suggest a pattern of learning that the model directly characterizes. With high statistical certainty, words require on the order of ~ 10 learning instances, which occur on average once every two months. Our method is extremely simple, statistically principled, and broadly applicable to modeling data-driven learning. The model is also able to predict the rate at which new words are learned, based on the amount of data required to learn them. This provides a precise way in which data results reveal that children are able to learn words in a highly efficient manner, despite the inherent ambiguity of language. We have collected a large cross-linguistic dataset to quantify the role of data in learning, and to compare the model both fits and predictions to other analyses of model performance. The parameters of the model are able to predict the rate at which new words are learned, based on the amount of data required to learn them. This provides a precise way in which data results reveal that children are able to learn words in a highly efficient manner, despite the inherent ambiguity of language. We have collected a large cross-linguistic dataset to quantify the role of data in learning, and to compare the model both fits and predictions to other analyses of model performance.

Item response theory

- A model of each child:
 - Has a latent parameter on the trait (θ)
 - Estimate of language ability
- A model of each word:
 - How difficult it is
 - Optionally: How well it discriminates low vs. high ability kids
- Allows prediction about new words and new children

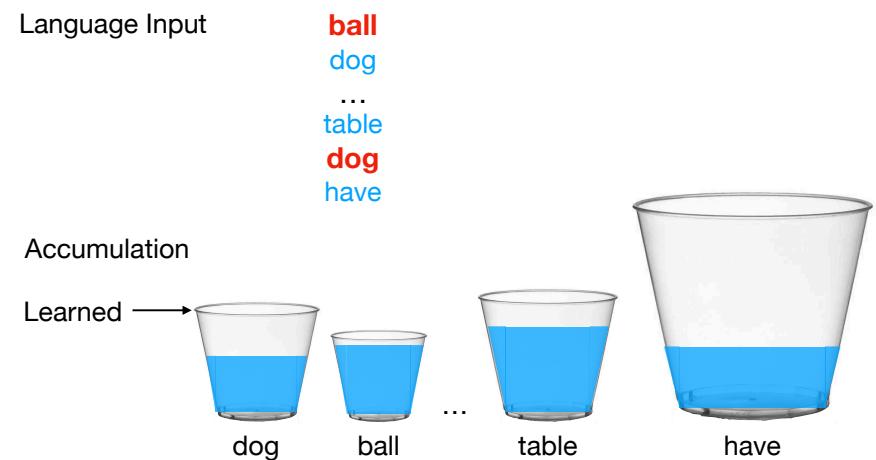


ICC = item characteristic curve

Embretson & Reise (2001)

Accumulator models are IRT models

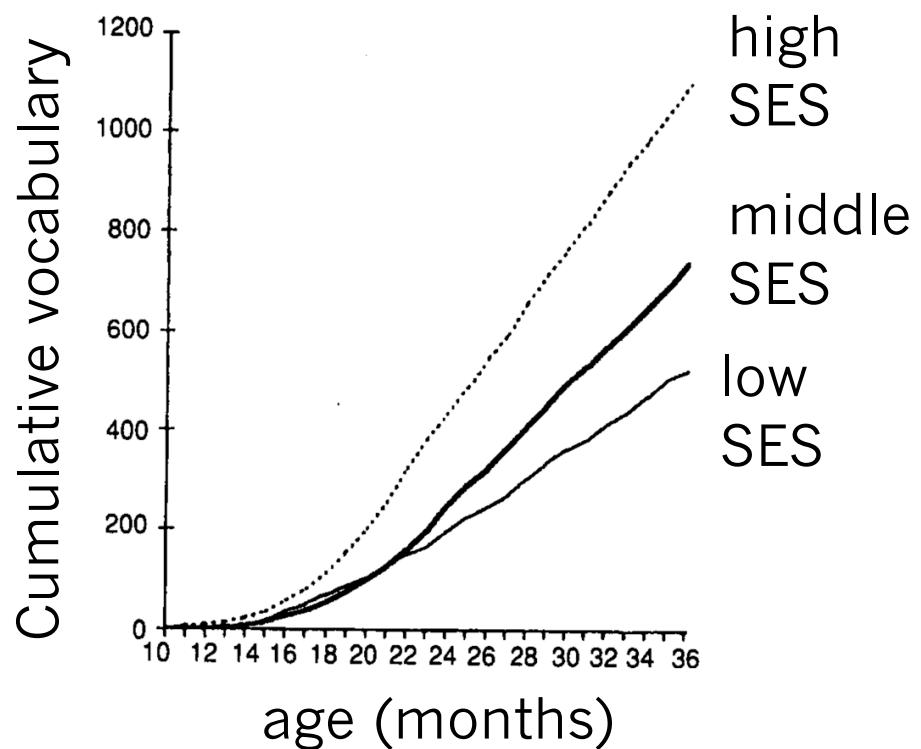
- Parameters map directly
 - Bucket size is word difficulty
 - Drip rate is child's learning rate X rate of input
- These models can be fit with mixed effects regression (de Boeck et al., 2011)
- But to be interpretable they require real measurements in absolute units
 - How many words per hour are heard?
 - How many total words are known?



Kachergis, Marchman, & Frank (2022)

Models are sometimes useful because they are wrong

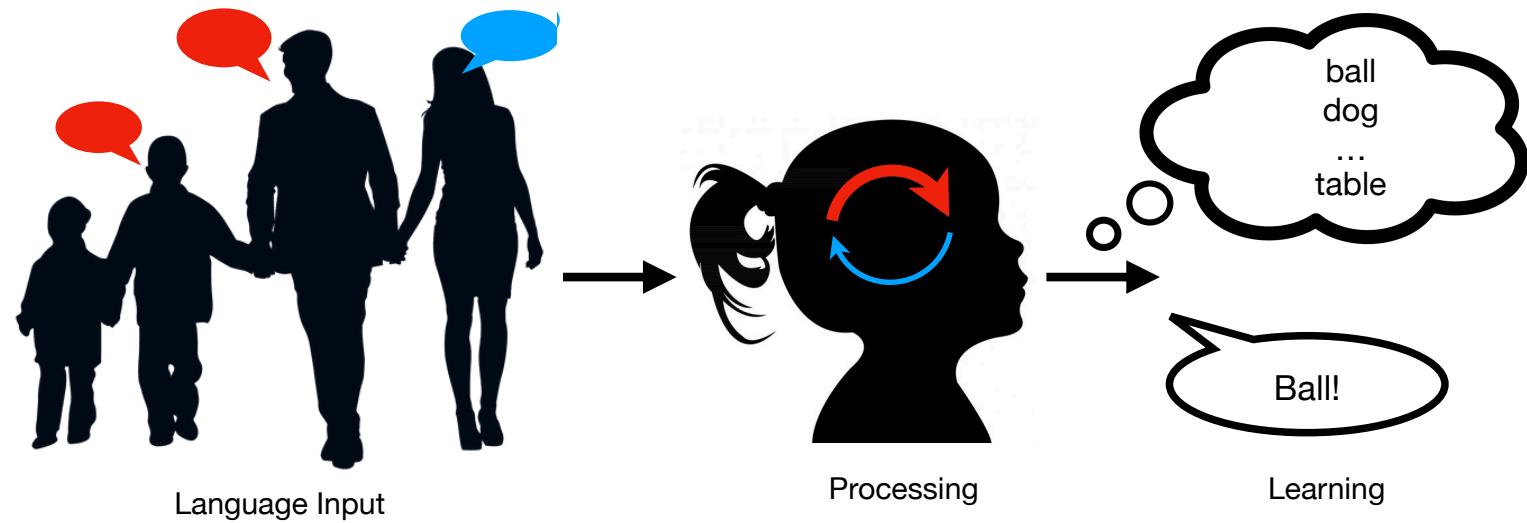
- Does input variation predict learning?
 - within culture, across populations?
 - across cultures?
- Failures to predict would constitute evidence for
 - processes beyond simple accumulation (e.g., learning from overhearing)
 - variation in learning process across populations



Hart & Risley (1995),
Cf. Raz & Beatty (2018), Sperry et al. (2018)

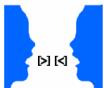
Outline

1. Basics of early word learning
2. Toward a “standard model” of language learning
- 3. Developing resources (to measure quantities in the model)**
 - Input – childe-db
 - Learning – wordbank
 - Processing – peekbank
4. The course toolset



Peekbank

CHILDES



Child Language Data Exchange
System

CHILDES is the child language component of the [TalkBank](#) system.
TalkBank is a system for sharing and studying conversational interactions.

System

[**Ground Rules**](#)
[Contributing New Data](#)
[IRB Principles](#)
[Overviews and Introductions](#)

Database

[**Index to Corpora**](#)
[Browsable Database](#)
[LuCiD Toolkit](#)
[childe-db](#)

Manuals

[CHAT - CLAN - MOR](#)
[Tutorial Screencasts](#)
[SLP's Guide to CLAN](#) and [中文](#)

The original “big data” for child language

[Other Child Language sites](#)

[CLAN](#)

Brian MacWhinney : [homepage](#)

... in a shared research environment with open tools and resources

[Unicode and IPA for Mac](#)

[Unicode and IPA for Windows](#)

Special Procedures

[CA analysis](#)

[Digitized video](#)

[Digitized audio](#)

Teaching Resources

[YouTube Examples](#)

[Bibliographies](#)

Versions

[Derived Corpora and Counts](#)

[XML version of the database](#)

[Database Versioning](#)

[MRC lexical dictionary](#)

ChildFREQ [Site](#) and [Paper](#)

More Resources

[Building a New Corpus](#)

[CCT Computerized
Comprehension](#)

[LEAT Assessment Tool](#)



childe-db

A flexible and reproducible interface to CHILDES



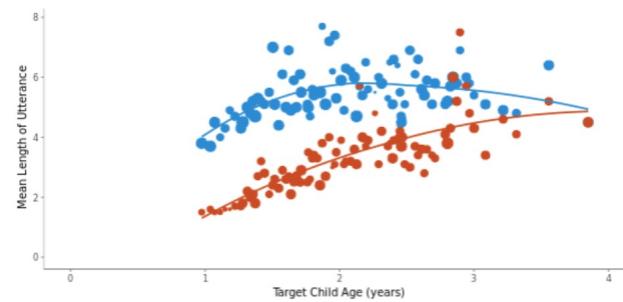
API Tutorial

Get a hands on walk-through on accessing [childe-db](#) through R.

```
> library(childebr)
> d_adam_prod <- get_tokens(collection = NULL,
+                               corpus = "Brown",
+                               role = "target_child",
+                               age = NULL,
+                               sex = NULL,
+                               child = "Adam",
+                               token = c("dog", "ball"))
Getting data from 1 child in 1 corpus ...
```

Visualizations

Explore the data in [childe-db](#) using our interactive applications.



Sanchez*, Meylan* et al. (2016), *Behavior Research Methods*

Accessing childesr-db

2019-10-16

Source: vignettes/access_childesr_db.Rmd

Overview

The `childesr` package allows you to access data in the childesr-db from R. This removes the need to write complex SQL queries in order to get the information you want from the database. This vignette shows some examples of how to use the data loading functions and what the resulting data look like.

There are several different `get_` functions that you can use to extract different types of data from the childesr-db:

- `get_transcripts()`
- `get_participants()`
- `get_tokens()`
- `get_types()`
- `get_utterances()`
- `get_speaker_statistics()`

Technical note 1: You do not have to explicitly establish a connection to the childesr-db since the `childesr` functions will manage these connections. But if you would like to establish your own connection, you can do so with `connect_to_childesr()` and pass it as an argument to any of the `get_` functions. If you do so, make sure to disconnect the connections you make by using `DBI::dbDisconnect()`, `childesr::clear_connections()`, or restarting your R session.

Technical note 2: We have tried to optimize the time it takes to get data from the database. But if you try to query and get all of the tokens, it will take a long time.

```
# load the library
library(childesr)
```

R API Overview

- Programmatic access to data
- Small number of calls to maintain
- Versioned database
- Caches lots of common data preprocessing
- Removes reproducibility uncertainty

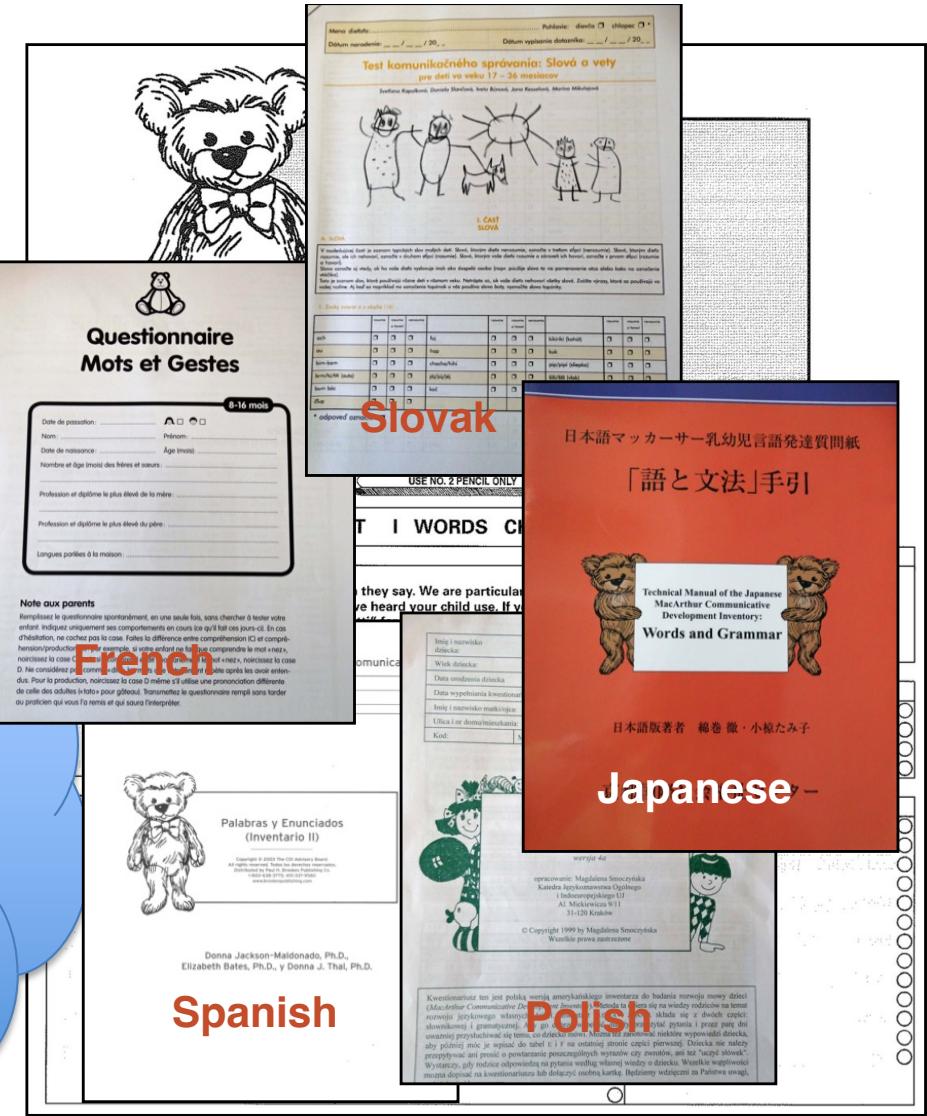
The MacArthur-Bates Communicative Development Inventory (CDI)

Ba,
da,
ma



Madeline (12 mo)

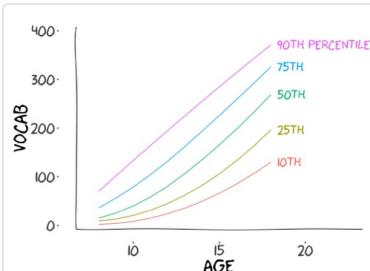
Hi, bye, ball,
dog, more,
milk,
mommy,
daddy, hand,
foot, nose, ...





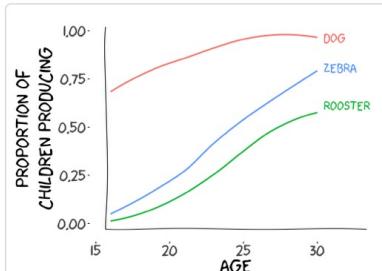
Wordbank

An open database of children's vocabulary development



Vocabulary Norms

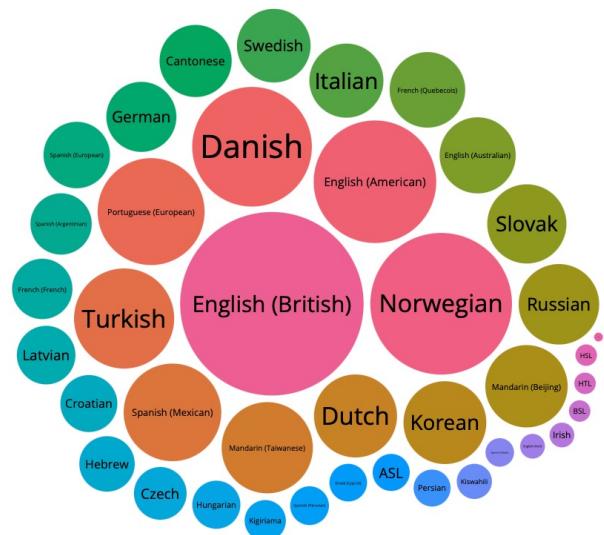
Explore vocabulary size growth curves for various languages and demographic groups.



Item Trajectories

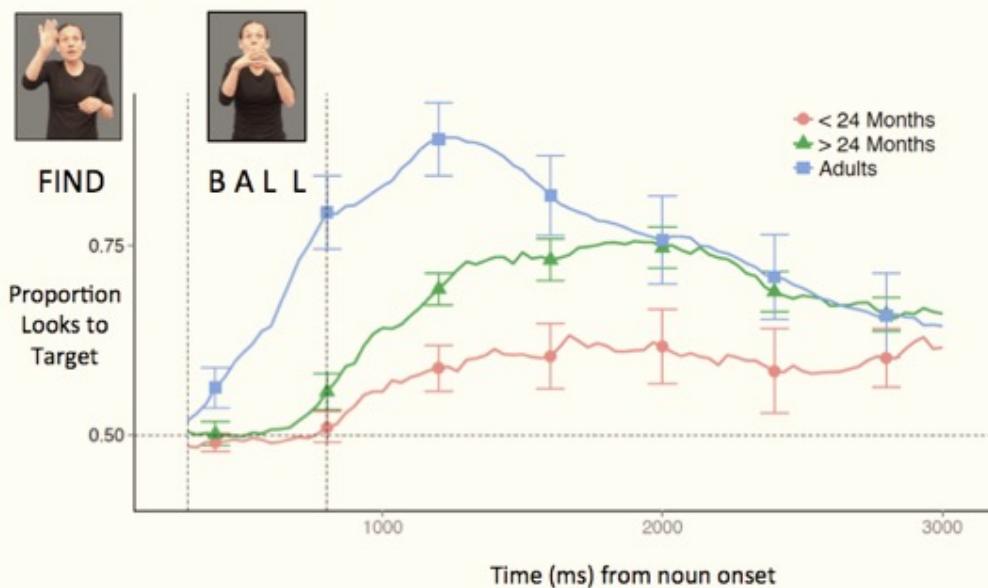
Explore trajectories of individual words, word categories, and grammar items.

Wordbank contains data from 79,027 children and 93,258 CDI administrations, across 40 languages and 76 instruments:



Frank et al. (2017), *Journal of Child Language*

Looking-While-Listening (LWL)



Parallel results in studies with diverse populations of children in the U.S.

- Spanish & English speakers in lower and higher SES families
- Late talkers
- Children born premature
- Emerging & simultaneous bilinguals
- Deaf children learning ASL

Slide courtesy Anne Fernald

MacDonald, Lamarr, Corina, Marchman, & Fernald (2018)



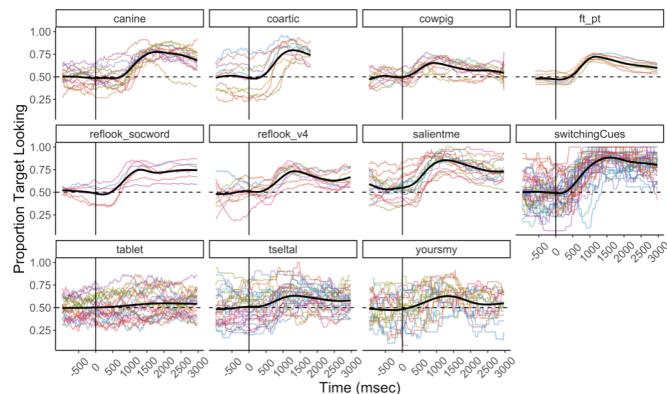
peekbank

A flexible and reproducible interface to developmental eyetracking datasets

What is peekbank?

peekbank is a flexible and reproducible interface to developmental eyetracking datasets.

The Peekbank project is an open database storing eye-tracking datasets on children's word recognition in a well-documented, easily accessible, tabular format. It also provides processing tools for standardizing eye-tracking data across data sources ([peekds R package](#)), interfaces for accessing the database ([peekbankr R package](#)), and applications for visualizing the data ([Peekbank Shiny App](#)).



Data Access Tutorial

Get started accessing the data here



Interactive Visualizations

Dynamically generate visualizations of the data



Documentation

Further information about the database

Outline

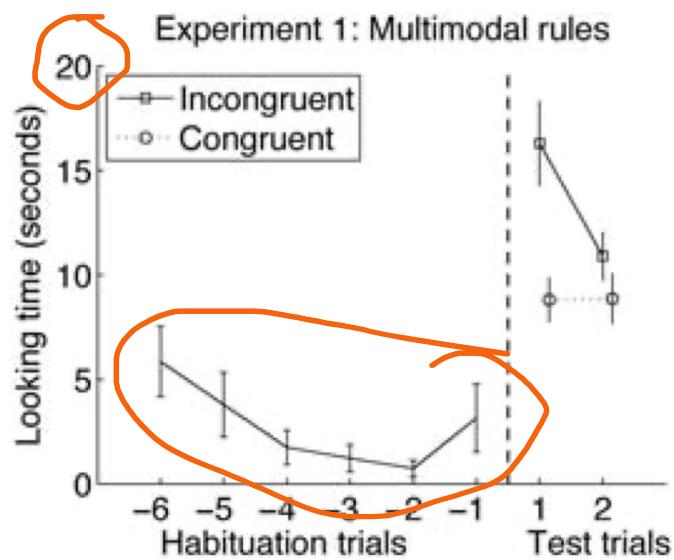
1. Basics of early word learning
2. Toward a “standard model” of language learning
3. Developing resources (to measure quantities in the model)
 - Input – childe-db
 - Learning – wordbank
 - Processing – peekbank
- 4. The course toolset**

		DATA	
		Same	Different
CODE	Same	Reproducible	Replicable
	Different	Robust	Generalizable

PAPER

Information from multiple modalities helps 5-month-olds learn abstract rules

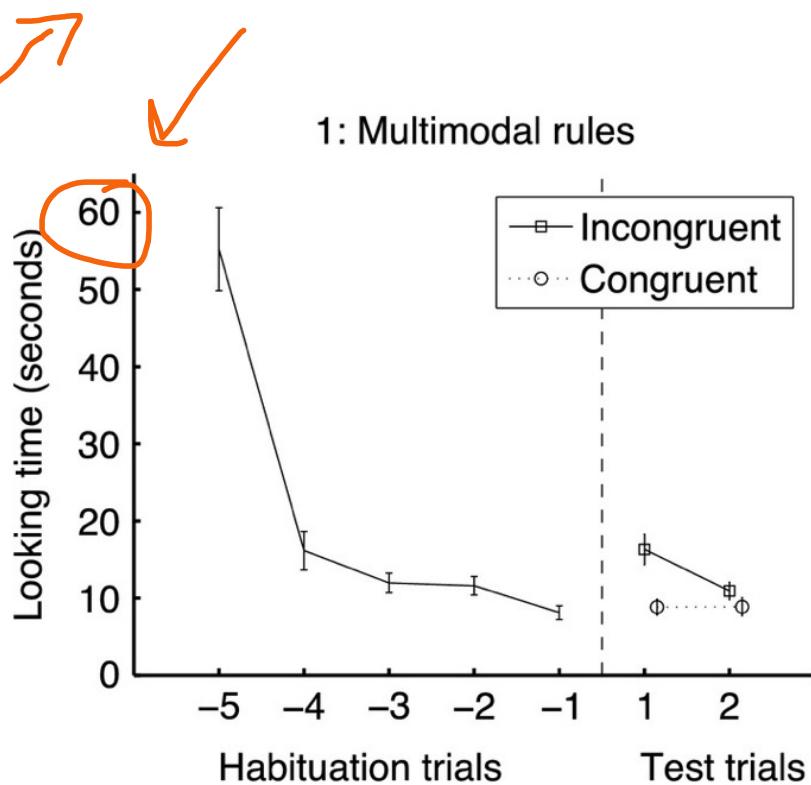
Michael C. Frank,¹ Jonathan A. Slemmer,² Gary F. Marcus³ and Scott P. Johnson⁴



Erratum

This article corrects the following: ▾

First published: 23 February 2013 | <https://doi.org/10.1111/desc.12060> | Citations: 1



Workflow

1. Get Excel data
2. Write Matlab code
3. Export figure as .EPS
4. Tweak in Illustrator
5. Drag into Word

```
%% new plots
expts = {'1: Multimodal rules','2: Visual rules','3: Auditory rules
(checkerboard)', '3: Auditory rules'};
expt_nums = [1 2 4];
s = {[1 2],[3 4],[6 7]};
figure(1); clf;
set(gcf,'Position',[ 121     48     885    694]);
```

Wrong trial numbers!

```
for i = 1:3
    subplot(2,4,s{i})
    set(gca,'FontSize',18);
    e = expt_nums(i);
    habit_times = [mean(HM7(EXP==e)) mean(HM8(EXP==e)) mean(HM9(EXP==e)) ...
        mean(HM10(EXP==e)) mean(HM11(EXP==e)) mean(HM12(EXP==e)))./1000;
    habit_errs = [stderr(HM7(EXP==e)) stderr(HM8(EXP==e)) stderr(HM9(EXP==e)) ...
        stderr(HM10(EXP==e)) stderr(HM11(EXP==e)) stderr(HM12(EXP==e)))./1000;
    hold on
    errorbar(1:6,habit_times(1:6),habit_errs(1:6),'-k')
    a = errorbar(7:8,[mean(N1C(EXP==e)) mean(N2C(EXP==e)))./1000, ...
        [stderr(N1C(EXP==e)) stderr(N2C(EXP==e)))./1000,'-sk');
    b = errorbar([7:8]+.15,[mean(F1C(EXP==e)) mean(F2C(EXP==e)))./1000, ...
        [stderr(F1C(EXP==e)) stderr(F2C(EXP==e)))./1000,:ok');
    axis([0 9 0 20]);
    if i == 1
        ylabel('Looking time (seconds)')
        xlabel('Habituation trials           Test trials','Position',[5 -2.2])
        legend([a,b],{'Incongruent','Congruent'},'Location','NorthWest')
    end
    title(expts{e});
    set(gca,'XTick',[1 2 3 4 5 6 7 8],'XTickLabel',{'-6','-5','-4','-3','-2','-
1','1','2'});
    line([6.5 6.5],[0 20],'LineStyle','--','Color',[0 0 0],'LineWidth',1)
end
```



Not just me...

```
gen recall11=.  
replace recall11=0 if Q21==1 replace recall11=1 if Q21==3 | Q21==5 | Q21==6  
replace recall11=2 if Q21==2 | Q21==4 | Q21==7 | Q21==8  
replace recall11=0 if Q69==1  
replace recall11=1 if Q69==3 | Q69==5 | Q69==6  
replace recall11=2 if Q69==2 | Q69==4 | Q69==7 | Q69==8 ta recall11
```

Corrigendum: Healthy Out-Group Members Are Represented Psychologically as Infected In-Group Members

First Published November 15, 2019 | Correction | [Find in PubMed](#) |  Check for updates
<https://doi.org/10.1177/0956797619887750>

[Article information](#) ▾



Original Article: [Healthy Out-Group Members Are Represented Psychologically as Infected In-Group Members](#)

Original article: Petersen, M. B. (2017). Healthy out-group members are represented psychologically as infected in-group members. *Psychological Science*, 28, 1857–1863. doi:[10.1177/0956797617728270](https://doi.org/10.1177/0956797617728270)

Andrew Lampinen, Kengthsagn Louis, and Michael Frank identified a coding error in the analyses reported in Petersen (2017). Specifically, when the summary measures of the number of within-group errors were generated in all three studies, a typo in the analysis code caused up to a third of the within-group errors for East Indian targets to be excluded. As a consequence, the number of within-group recall errors was underreported in the experimental condition for all three studies. This Corrigendum will correct those errors.

Low data and code availability renders most papers non-reproducible

Data & code availability (prevalence estimates)

	Data	Code
Psychology (2014-2017) ¹	2% [1-4%]*	1% [0-1%]
Social Sciences (2014-2017) ²	7% [2-13%]	1% [0-3%]
Biomedicine (2015-2017) ³	18.3% [11.6- 27.3%]	0%

Data not available on request

	Data shared
141 articles published in four major APA journals (2004) ⁴	27%
516 ecology articles published (1991-2011) ⁵	20%
111 most highly-cited articles published in psychology & psychiatry (2006-2016) ⁶	14%

¹Hardwicke et al. (2020a)

²Hardwicke et al. (2020b)

³Wallach et al. (2018)

*[95% confidence intervals]

⁴Wicherts et al. (2006)

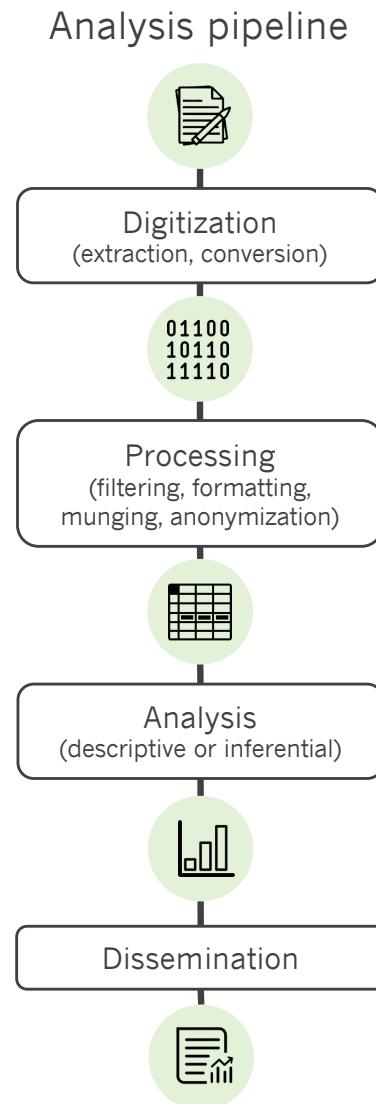
⁵Vines et al. (2014)

⁶Hardwicke & Ioannidis (2018)

What is a reproducible report?

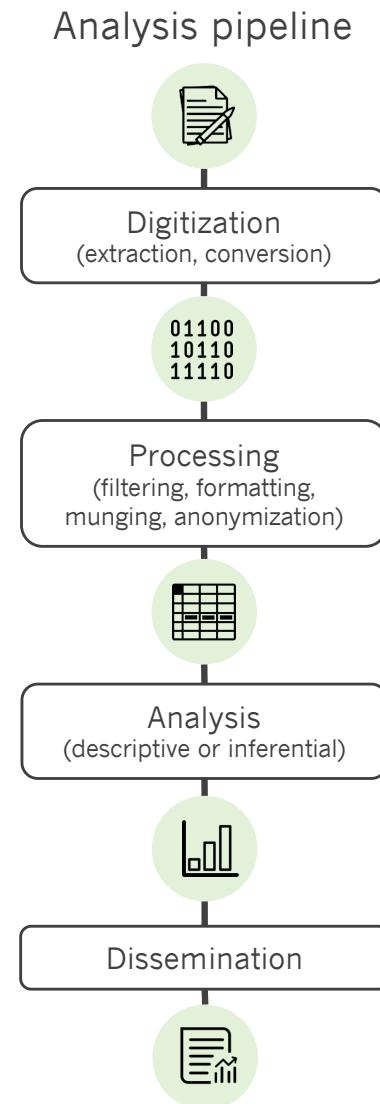
“...a scientific publication is not the scholarship itself, it is merely advertising of the scholarship. The actual scholarship is the complete software development environment and the complete set of instructions which generated the figures.”

- Buckheit & Donoho (1995)



What is a reproducible report?

A document for which the entire analysis pipeline, from primary data to reported numerical values, can be recreated, ideally by an independent scientist.



Why write reproducible reports?

1. Provenance, transparency, & verification

To know what numbers mean, we need to know where they are from.

Transparent sharing of analysis pipelines enables independent verification.

2. Error detection & mitigation

Scientists are human and humans make mistakes. Disorganized workflows with manual tasks, create opportunities for error.

In a reproducible report, the flow of information through the analysis pipeline is directed and documented by code. This is easier to audit and manual tasks like copy & pasting are avoided, reducing opportunities for error.

3. Efficiency

Scientists often need to re-run or re-use analyses, for example in responding to reviewers, repurposing for a new project, or generating reporting outputs (posters, slides, papers etc.)

In a reproducible workflow, it is easier to build on your prior work by re-using code – there's no need to start from scratch each time or remind yourself what you did previously.

Our course repository

The screenshot shows a GitHub repository page for 'mcfrank/lot-language-learning-2023'. The repository is public and has 7 commits. The main branch is 'main'. The repository contains files like 'data', '.gitignore', 'LICENSE', 'README.md', 'day-1-tidyverse.Rmd', and 'lot-language-learning-2023.Rproj'. The 'About' section describes the repository as 'Materials for LOT School 2023, "Language Learning: A Data-Driven Approach"'. It includes links to 'Readme', 'MIT license', '1 star', '1 watching', and '0 forks'. The 'Releases' section indicates 'No releases published' and 'Create a new release'. The 'Packages' section indicates 'No packages published' and 'Publish your first package'. The 'README.md' file content is as follows:

```
lot-language-learning-2023

Materials for LOT School 2023, "Language Learning: A Data-Driven Approach"

Instructor: Michael C. Frank (Stanford)

Course Description

In this course, we will examine early language learning through the lens of new data resources that facilitate quantitative studies. Our framework will be the "Standard Model" of Kachergis, Marchman, and Frank (2022) that links language input to processing and learning outcomes, and we will consider the strengths and weaknesses of this model for describing vocabulary learning as well as the learning of some morphology and syntax. Our hands-on approach will involve learning the use of CHILDES and childe-db for studying language input, Wordbank for
```

Git and github

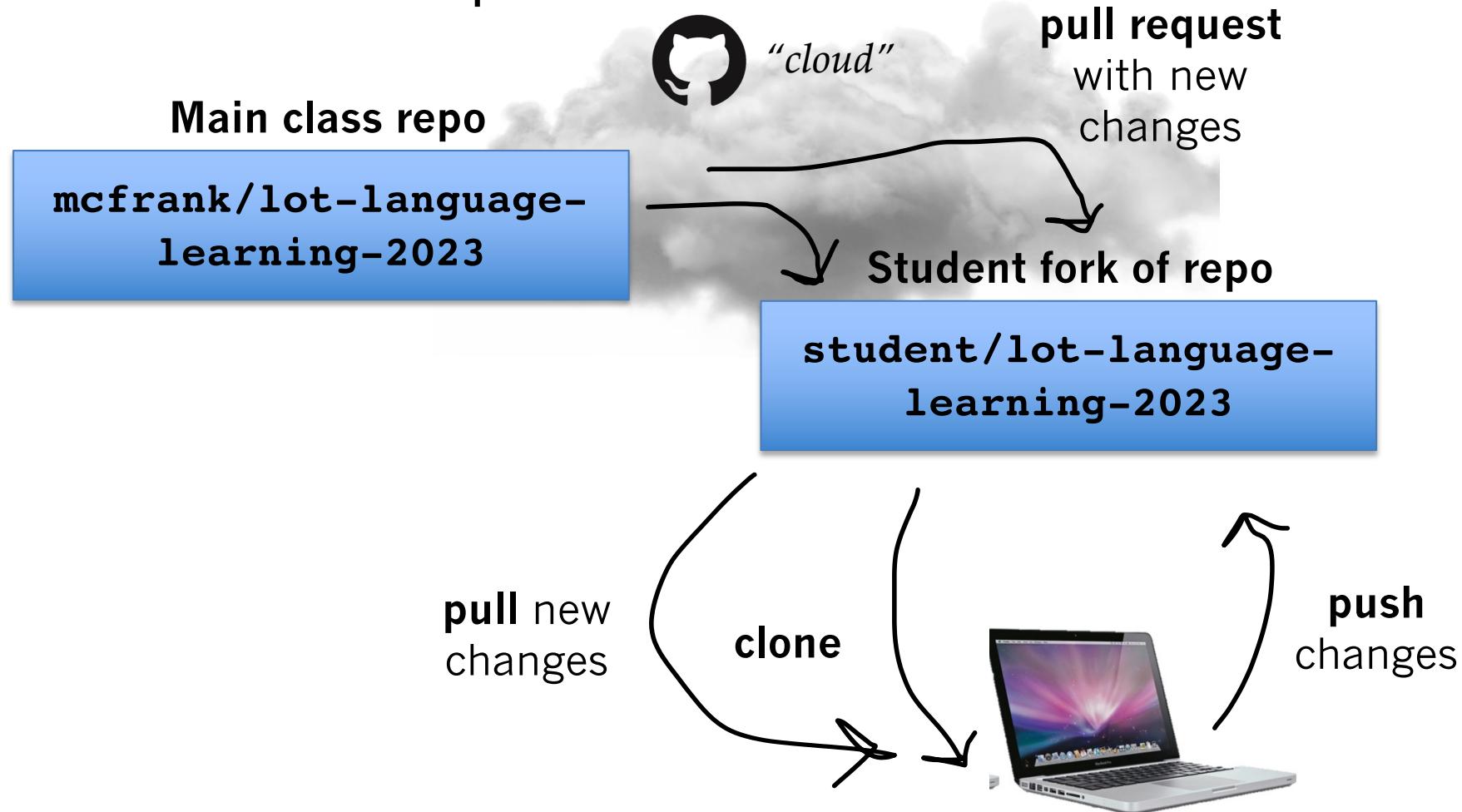
But, we don't want students
to alter the main
repository...



Two ways to engage with our course repo

1. Download the repo directly onto your computer
 - Each day, copy the next day's code into a new Rmd file
2. "Fork" the repo on github (a bit trickier, but helps you master a good workflow)
 - Clone your fork onto your own computer
 - When you are done for the day, commit your changes to your local repo
 - The next day, pull changes from my repo to your fork
 - Then pull the changes to your local computer

The “fork and pull” model



R Markdown

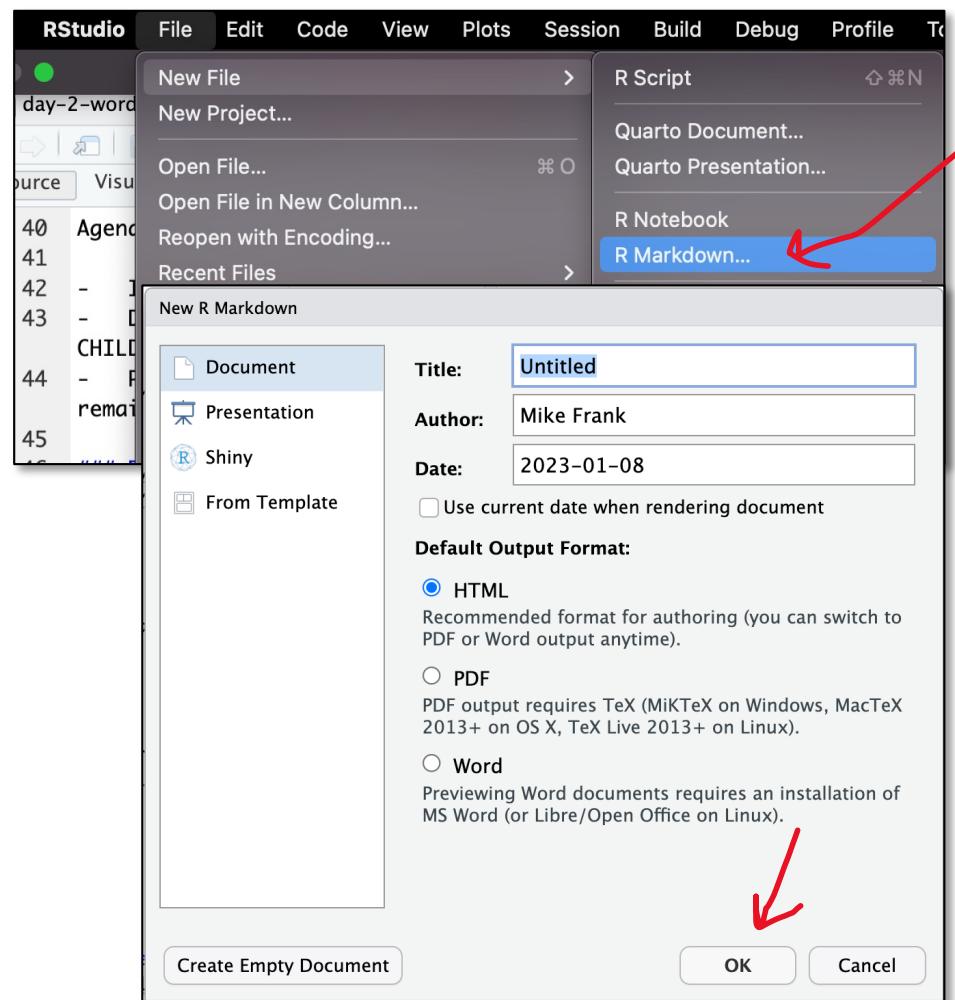
YAML header – don't mess with this

```
day-1-tidyverse.Rmd x
Source Visual
1 ---  
2 title: 'Day 1: Tidyverse refresher'  
3 author: "Mike Frank"  
4 date: "2023-01-09"  
5 output: html_document  
6 ---  
7  
8 ``{r}  
9 library(tidyverse)  
10 ``  
11  
12 This is a quick tidyverse introduction/refresher, adapted for the LOT Language Learning course  
from Appendix C of [Experimentology](http://experimentology.io). The topics it covers are:  
13  
14 - so-called "tidy" data  
15 - pipes (`%>%` and `|>`)  
16 - a few tidyverse verbs (`filter`, `mutate`, `summarise`, and `group_by`)  
17 - the barest bit of visualization using `ggplot`, and  
18 - joining tidy data frames.  
19  
20 I assume as a prerequisite that you already have some familiarity with R and am really just  
giving a refresher so everyone is on the same page. The best reference for this material is  
Hadley Wickham's [R for data scientists](http://r4ds.had.co.nz/) and I encourage you to read it
```

R code, will be executed when you hit this:

Formatted markdown text

You try it!



The screenshot shows the RStudio interface with the 'Source' tab selected. A red circle highlights the 'Knit' button in the toolbar. The code editor displays an R Markdown document with the following content:

```
1 ---  
2 title: "Untitled"  
3 author: "Mike Frank"  
4 date: "2023-01-08"  
5 output: html_document  
6 ---  
7  
8 ````{r setup, include=FALSE}  
9 knitr::opts_chunk$set(echo = TRUE)  
10 ````  
11  
12 ## R Markdown  
13  
14 This is an R Markdown document. Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents. For more details on using R Markdown see http://rmarkdown.rstudio.com.  
15  
16 When you click the **Knit** button a document will be generated that includes both content as well as the output of any embedded R code chunks within the document. You can embed an R code chunk like this:  
17  
18 ````{r cars}  
19 summary(cars)  
20 ````  
21  
22 ## Including Plots  
23  
24 You can also embed plots, for example:  
25  
26 ````{r pressure, echo=FALSE}  
27 plot(pressure)  
28 ````  
29  
30 Note that the `echo = FALSE` parameter was added to the code chunk to prevent printing of the R code that generated the plot.  
31
```

The tidyverse...

- Open `lot-language-learning-2023.Rproj`
 - Project file ensures that you will be working from the correct directory
- Open `day-1-tidyverse.Rmd`
- We will gradually work through each topic, pausing for simple exercises

Exercises

I have a problem,
please come help!

I'm done with this
exercise and ready to
move on