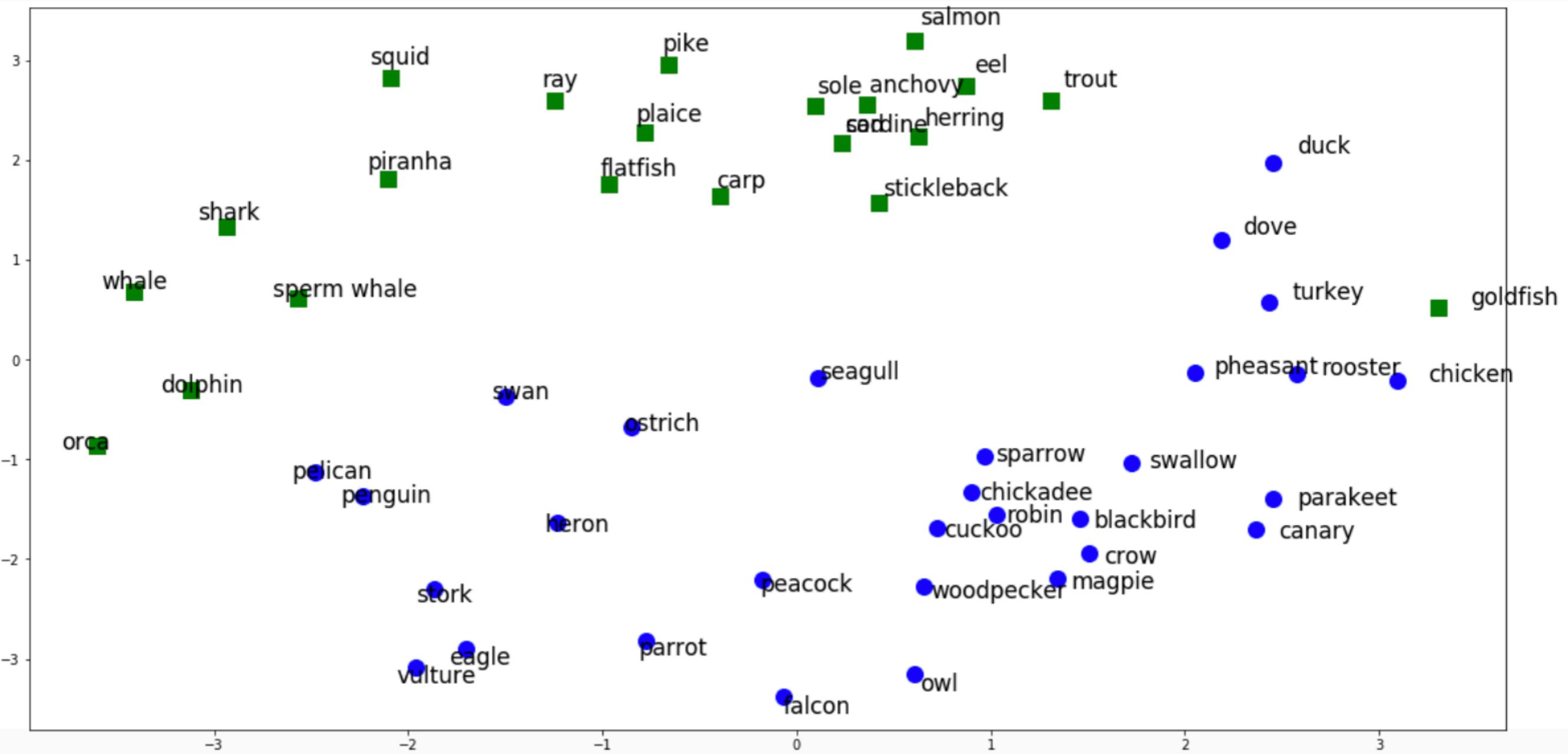


COG260: Data, Computation, and The Mind Languages

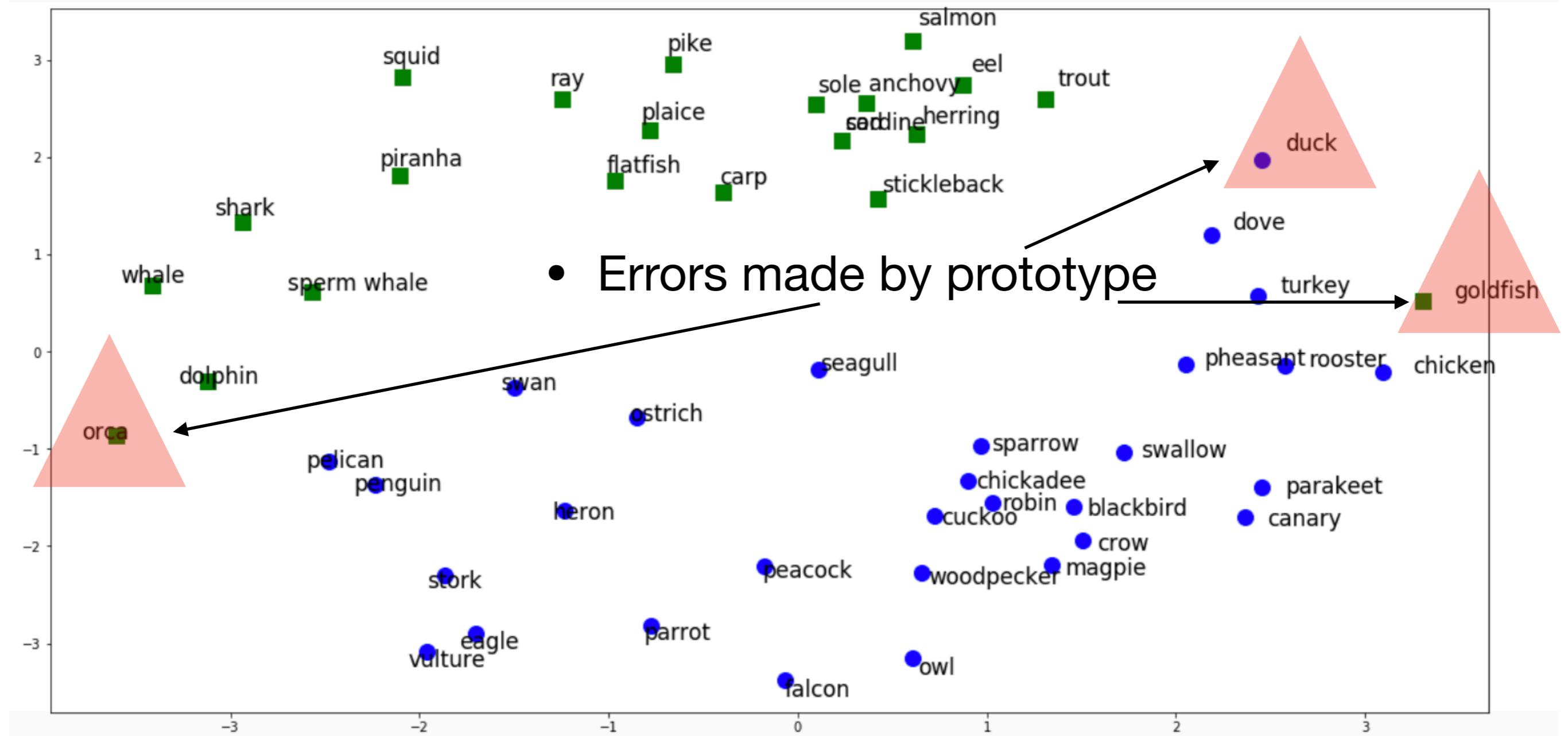
Lab 5: Categorization

- Prototype model: ~94%
- Exemplar model: ~96%

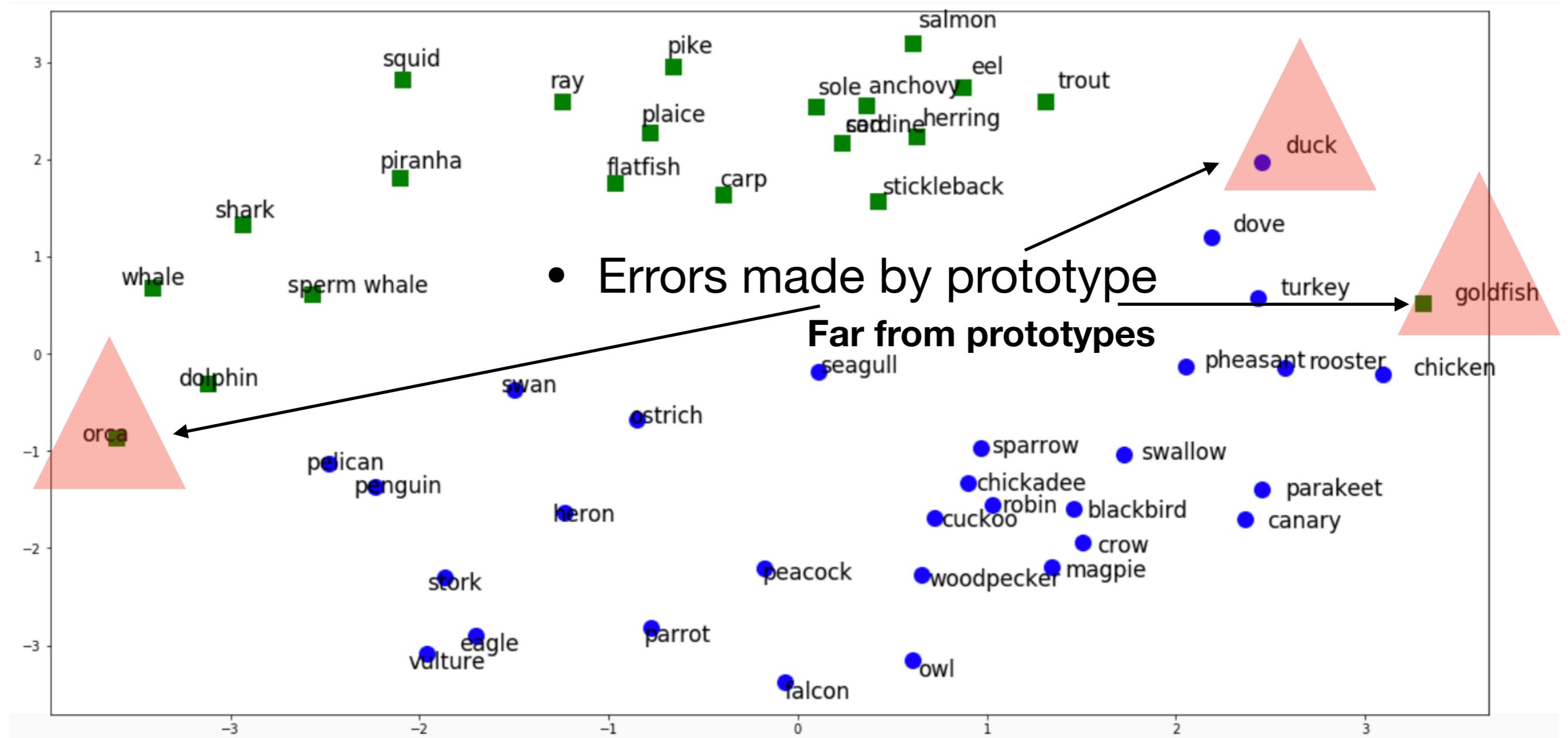
Lab 5: Categorization



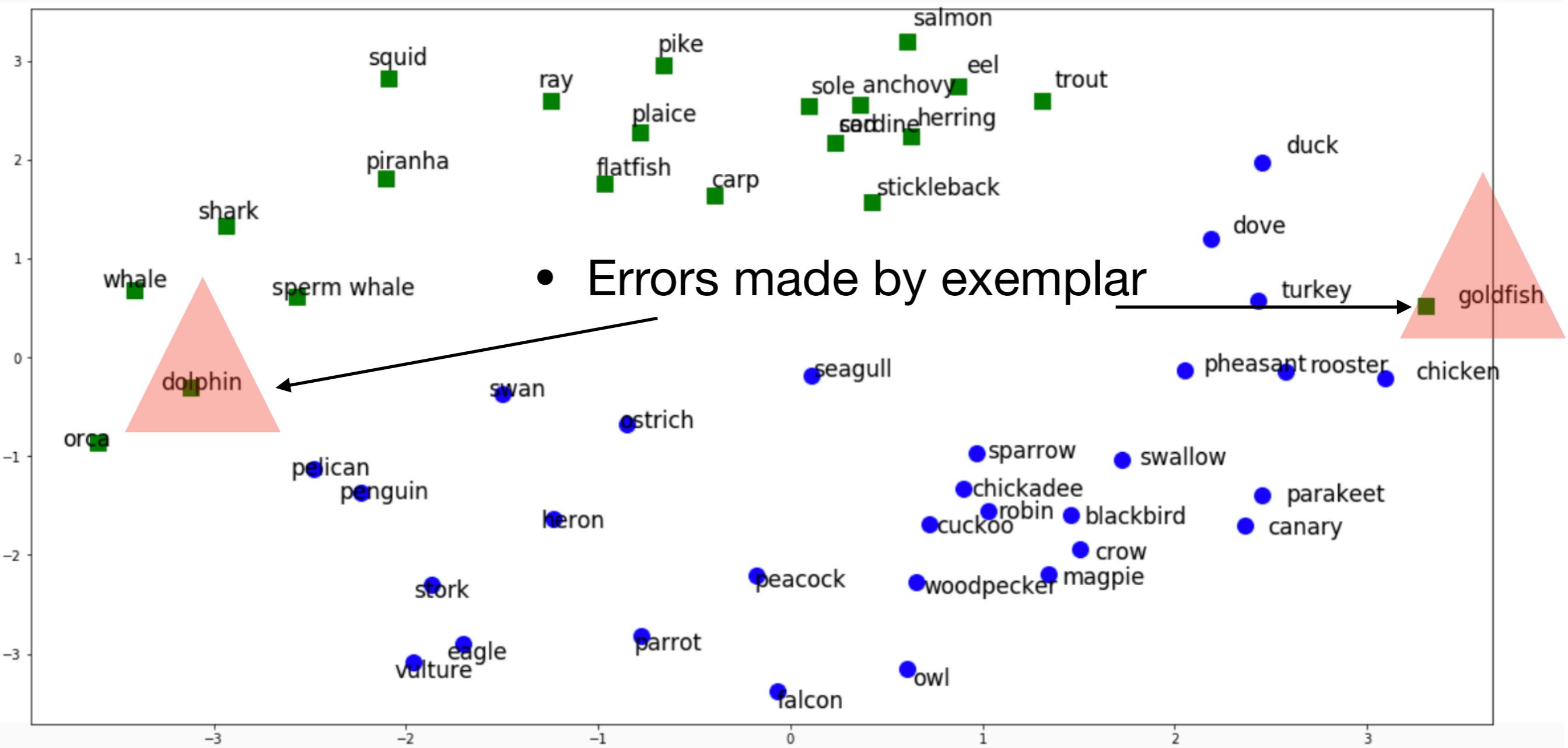
Discussion 1: What went wrong?



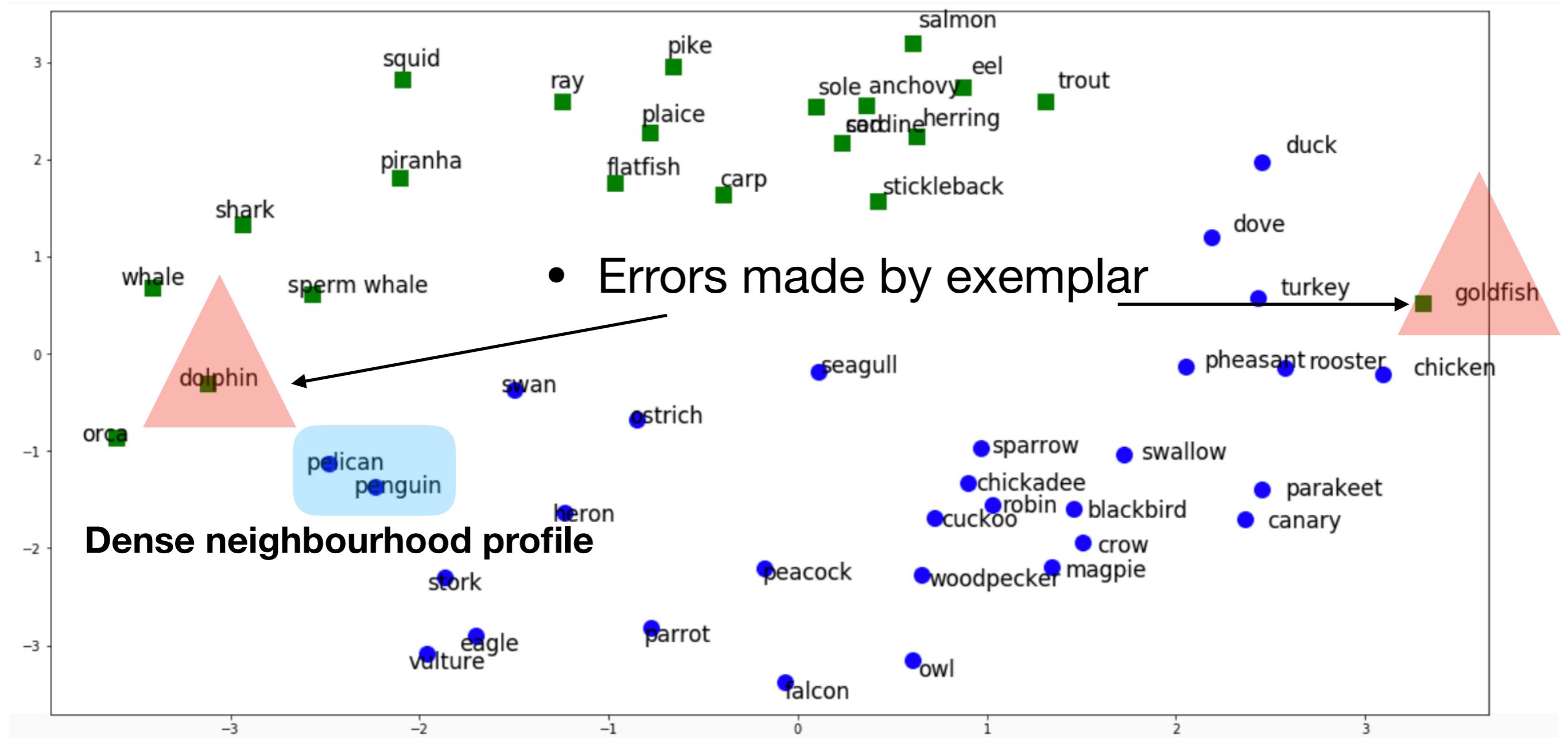
Misprediction of outliers



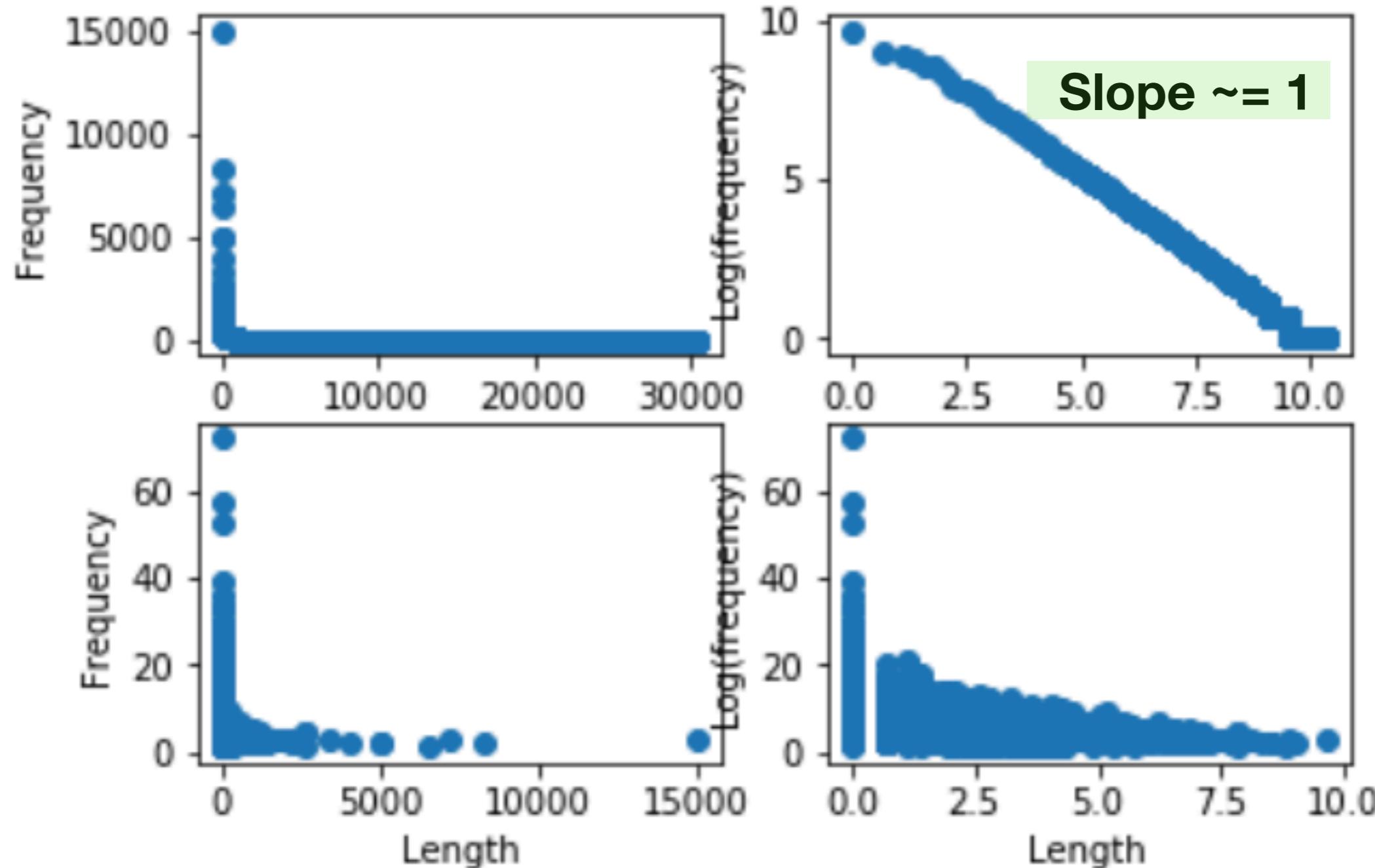
Discussion 2: What went wrong?



The lure of close neighbours

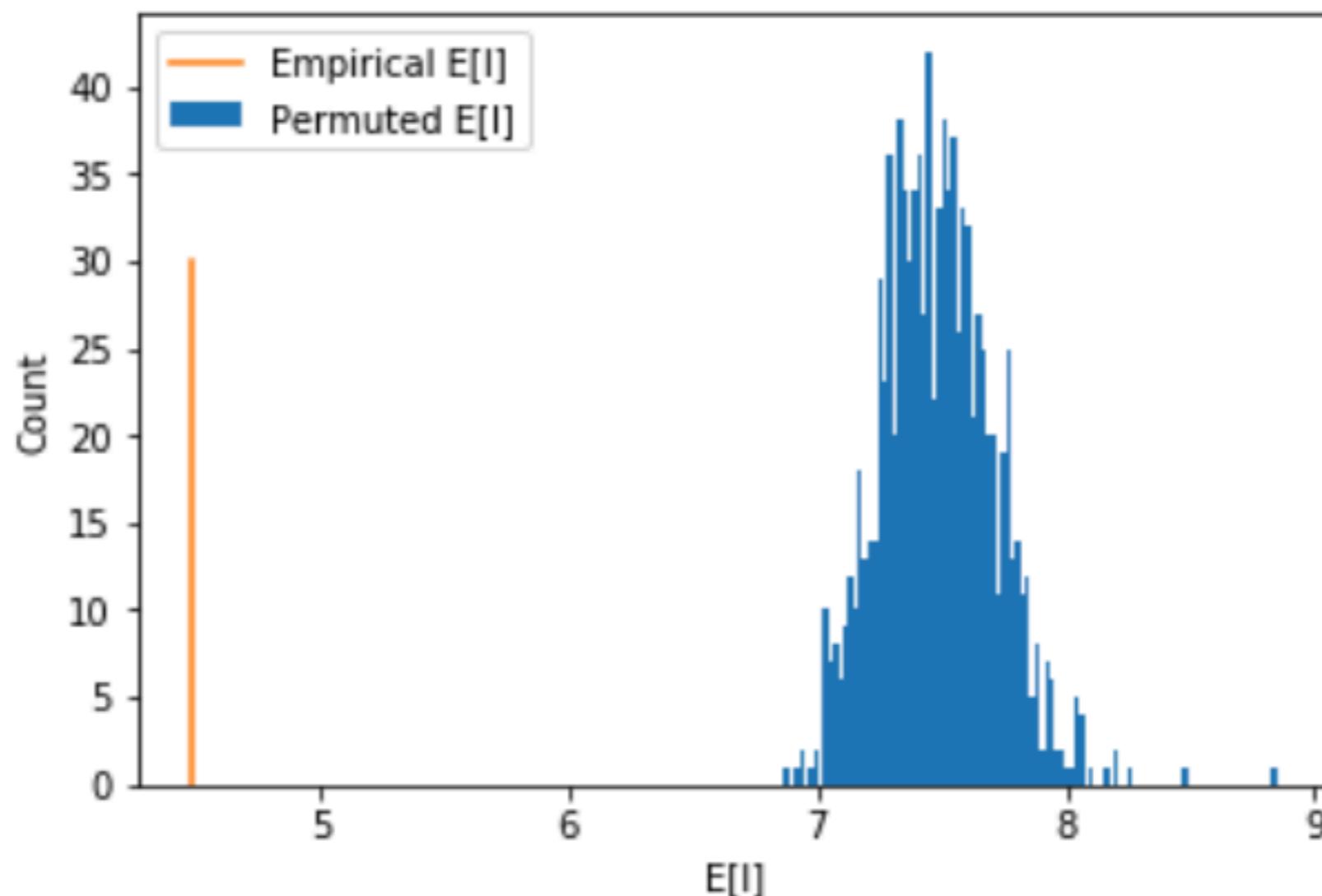


Lab 6: Word frequency



Lab 6: Word frequency

Expected length = 4.48533473678
p-value = 0.0



Suppose there are only 6 words in a hypothetical piece of text:

dog, dog, dog, wags, wags, tail.

1. How to calculate the “expected word length”:

By definition, the expected word length is equivalent to the average word length. Given the 6 words, we have the following word length distribution:

dog	dog	dog	wags	wags	tail
3	3	3	4	4	4

Therefore, the average is simply: $(3+3+3+4+4+4) / 6 = 3.5$

Alternatively, this quantity can also be calculated with the formula: $E[X] = \sum x \cdot p(x)$. Here x = length, and $p(x)$ = normalized frequency of a word type. Note in the 6-word text, there are only 3 unique word types with the following frequency distribution:

dog	wags	tail
3	2	1

Hence to calculate the expected word length, we can first calculate $p(x)$ or normalized frequency for each of the above word types:

dog	wags	tail
3/6	2/6	1/6

Finally, we can sum the product of these frequencies and word-type lengths to obtain the expected length:

$$E[\text{length}] = 3/6 \times 3 + 2/6 \times 4 + 1/6 \times 4 = (3 \times 3 + 2 \times 4 + 1 \times 4) / 6 = 3.5$$

2. How to obtain the “permuted expected word lengths”:

We can permute the frequency distribution of the three word types by exhaustively shuffling their frequencies in this case:

dog	wags	tail	
3	2	1	permutation 1
2	3	1	permutation 2
1	2	3	permutation 3
3	1	2	permutation 4
2	1	3	permutation 5
1	3	2	permutation 6

For each of these permuted trials, we can calculate the expected word length in the same way:

$$\text{permutation 1: } \frac{3}{6} \times 3 + \frac{2}{6} \times 4 + \frac{1}{6} \times 4 = \frac{(3 \times 3) + (2 \times 4) + (1 \times 4)}{6} = 3.5$$

$$\text{permutation 2: } \frac{2}{6} \times 3 + \frac{3}{6} \times 4 + \frac{1}{6} \times 4 = \frac{(2 \times 3) + (3 \times 4) + (1 \times 4)}{6} = 3.66666$$

$$\text{permutation 3: } \frac{1}{6} \times 3 + \frac{2}{6} \times 4 + \frac{3}{6} \times 4 = \frac{(1 \times 3) + (2 \times 4) + (3 \times 4)}{6} = 3.83333$$

permutation 4: $3/6 \times 3 + 1/6 \times 4 + 2/6 \times 4 = (3 \times 3 + 1 \times 4 + 2 \times 4) / 6 = 3.5$

permutation 5: $2/6 \times 3 + 1/6 \times 4 + 3/6 \times 4 = (2 \times 3 + 1 \times 4 + 3 \times 4) / 6 = 3.66666$

permutation 6: $1/6 \times 3 + 3/6 \times 4 + 2/6 \times 4 = (1 \times 3 + 3 \times 4 + 2 \times 4) / 6 = 3.83333$

3. How to calculate the *p*-value:

In this test, we make the proposal that the expected word length in the English lexicon is short, and hence we calculate the *p*-value by considering the proportion of permuted trials that have expected length SMALLER THAN OR EQUAL TO the actual expected length (i.e., we need to consider the extreme cases where the permuted trials produced even shorter expected length than the actual expected length). Because 3.5 showed up twice in the calculated expected word lengths in Section 2 and there was no trial that yielded a length shorter than 3.5, the *p*-value is then:

$$p = (\#\text{permuted length} \leq 3.5) / \#\text{permutations} = 2 / 6 = 1/3 = .33333$$

4. How to conclude from the *p*-value:

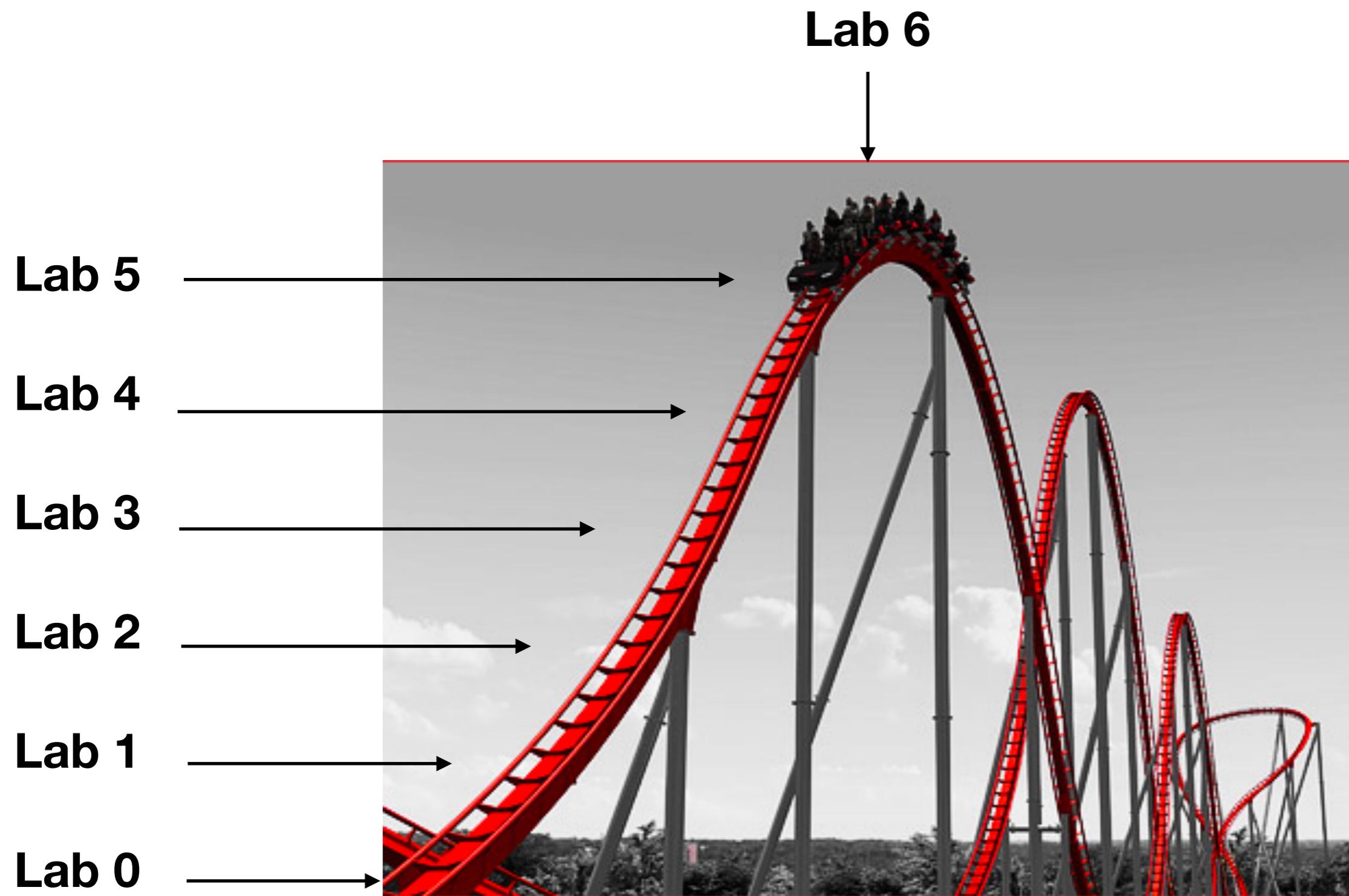
The null hypothesis of the permutation test is that there is no difference between the actual expected word length and those observed from the hypothetically constructed text (i.e. from frequency shuffles). Typically, we reject the null hypothesis at $p < 0.05$. Because the actual *p*-value > 0.05 , we cannot reject the null in this hypothetical case. In the actual lab, because the number of words is far greater than 6, you would not be able to exhaust the permutation trials, hence why you will only try for 1000 permutation trials. But the calculations of the expected word length and *p*-value follow the same procedures as above.

Lab 6: Word frequency

At least 1 of the permutations (even though it might not be within 1,000 samples) should yield a lower expected word length than the observed.

An alternative way of thinking about this problem:
What kind of word-frequency mapping would yield the lowest expected length?

Up to now



Up to now



Lab 0 → Piechart

Up to now



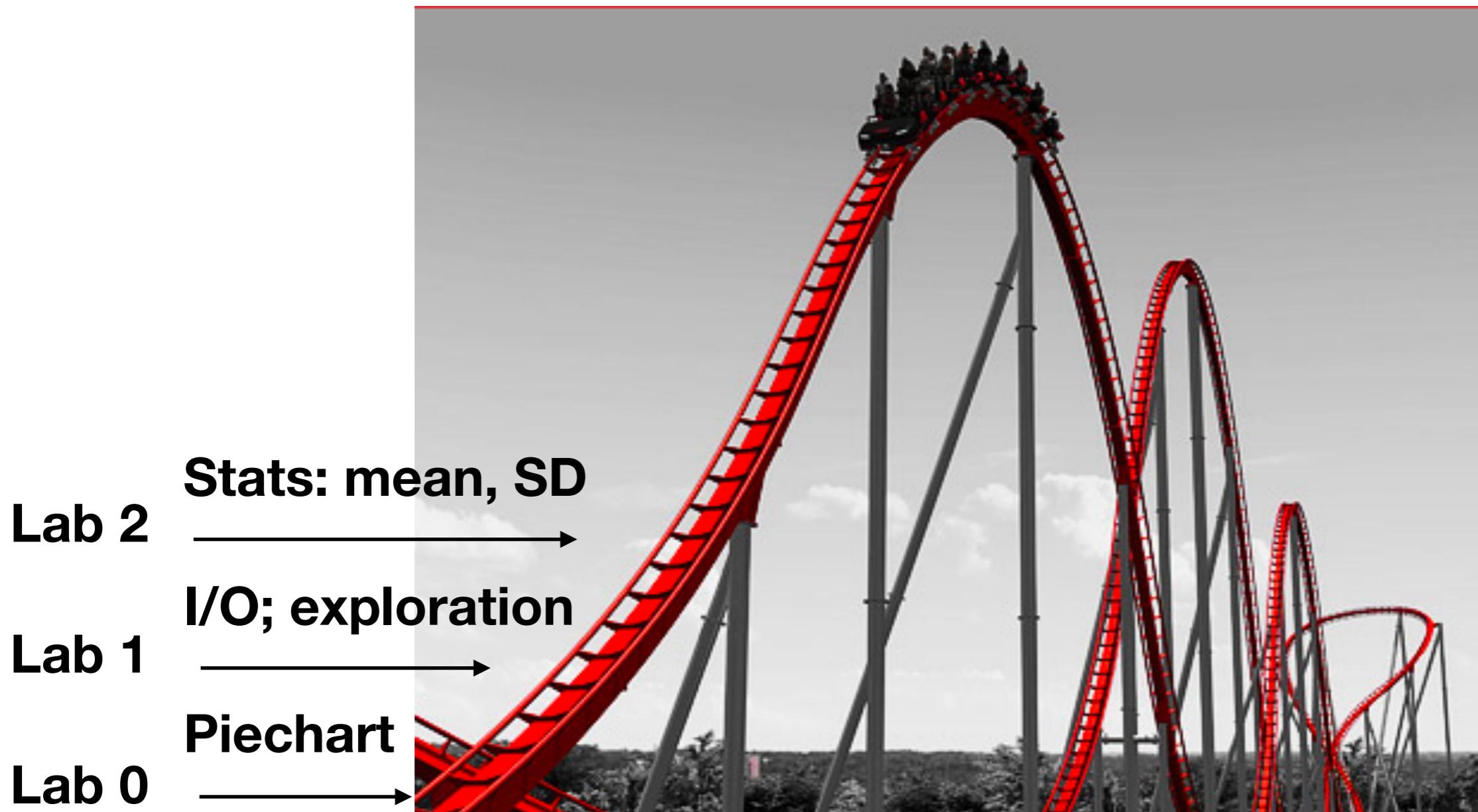
Lab 1

I/O; exploration

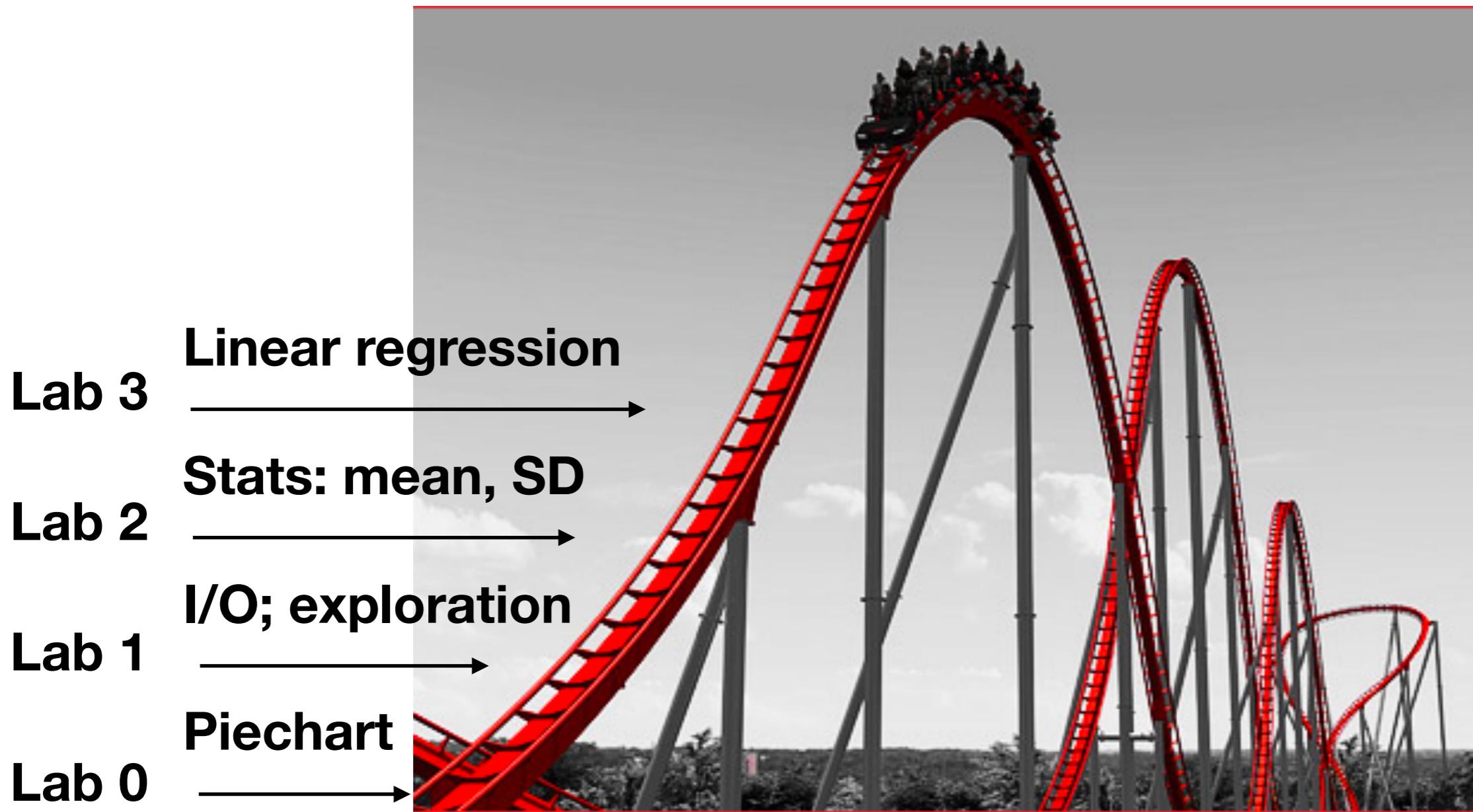
Lab 0

Piechart

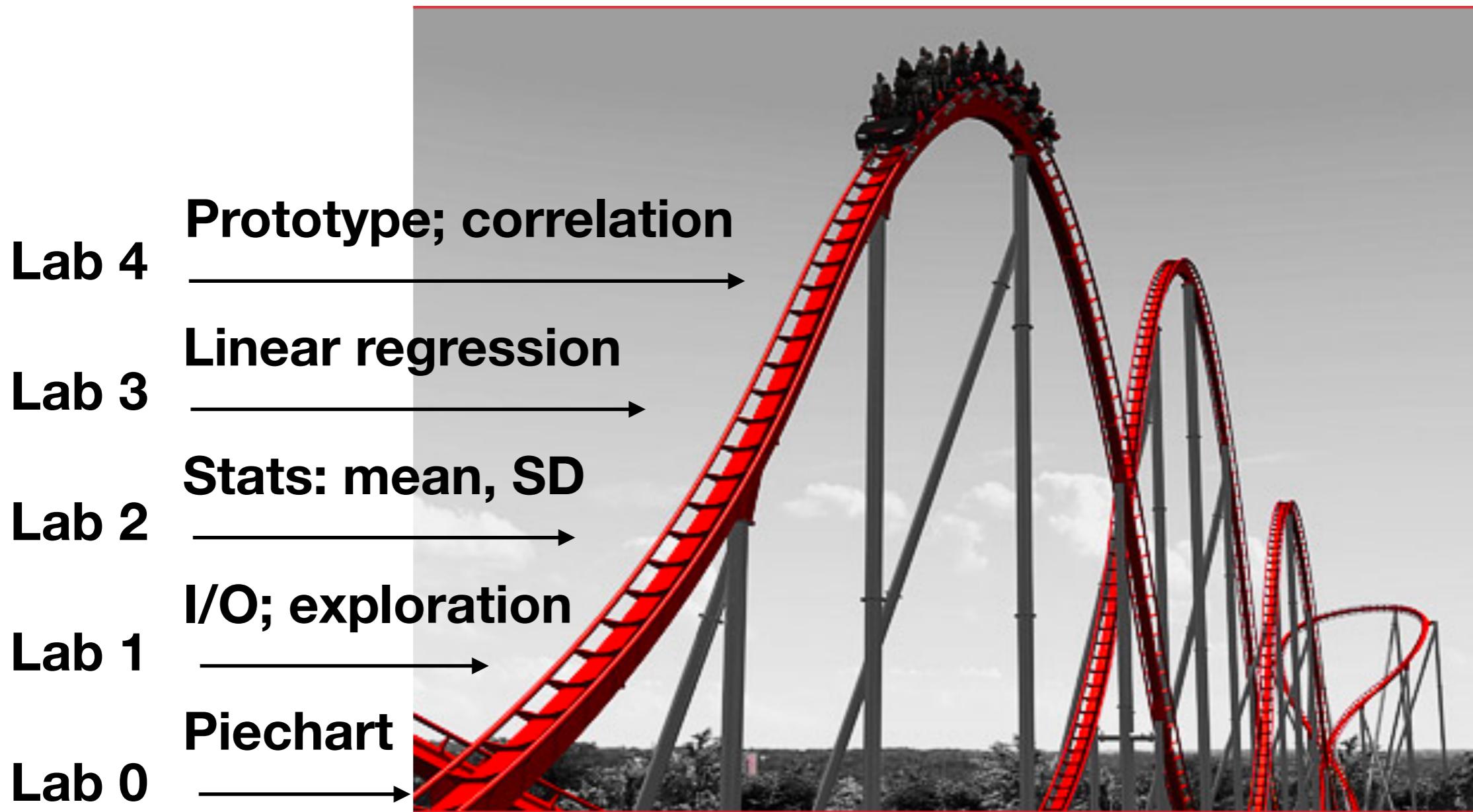
Up to now



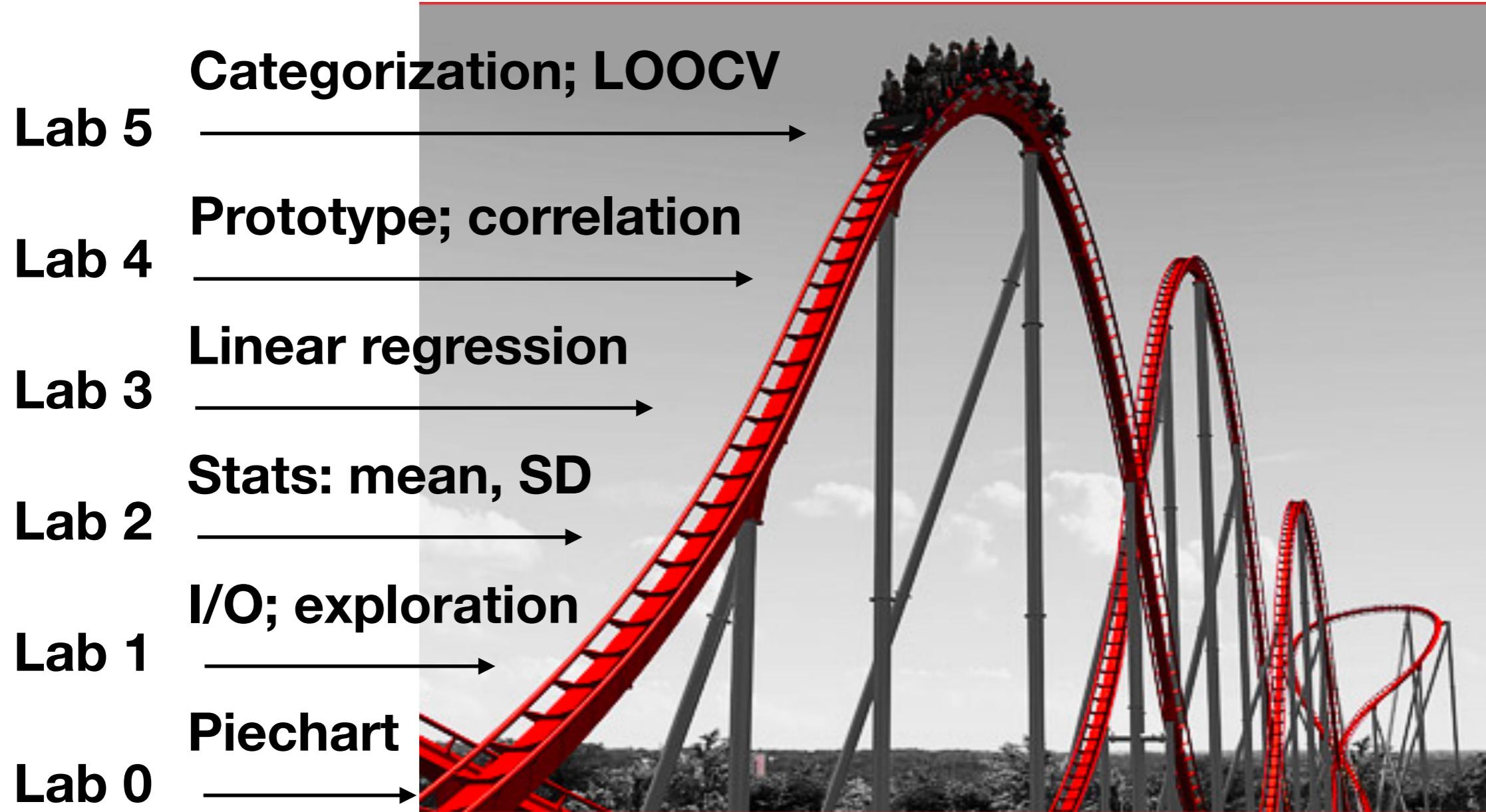
Up to now



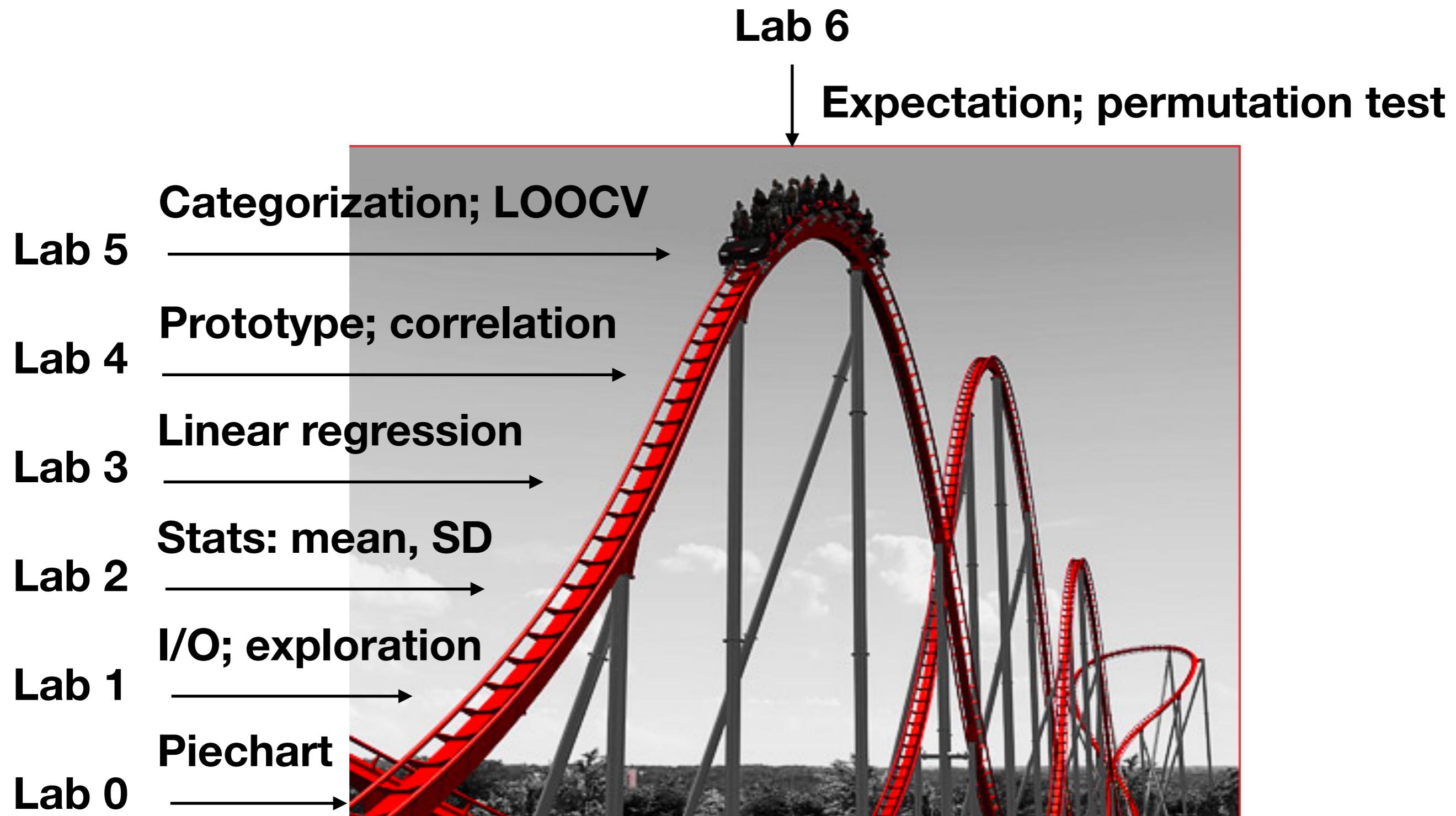
Up to now



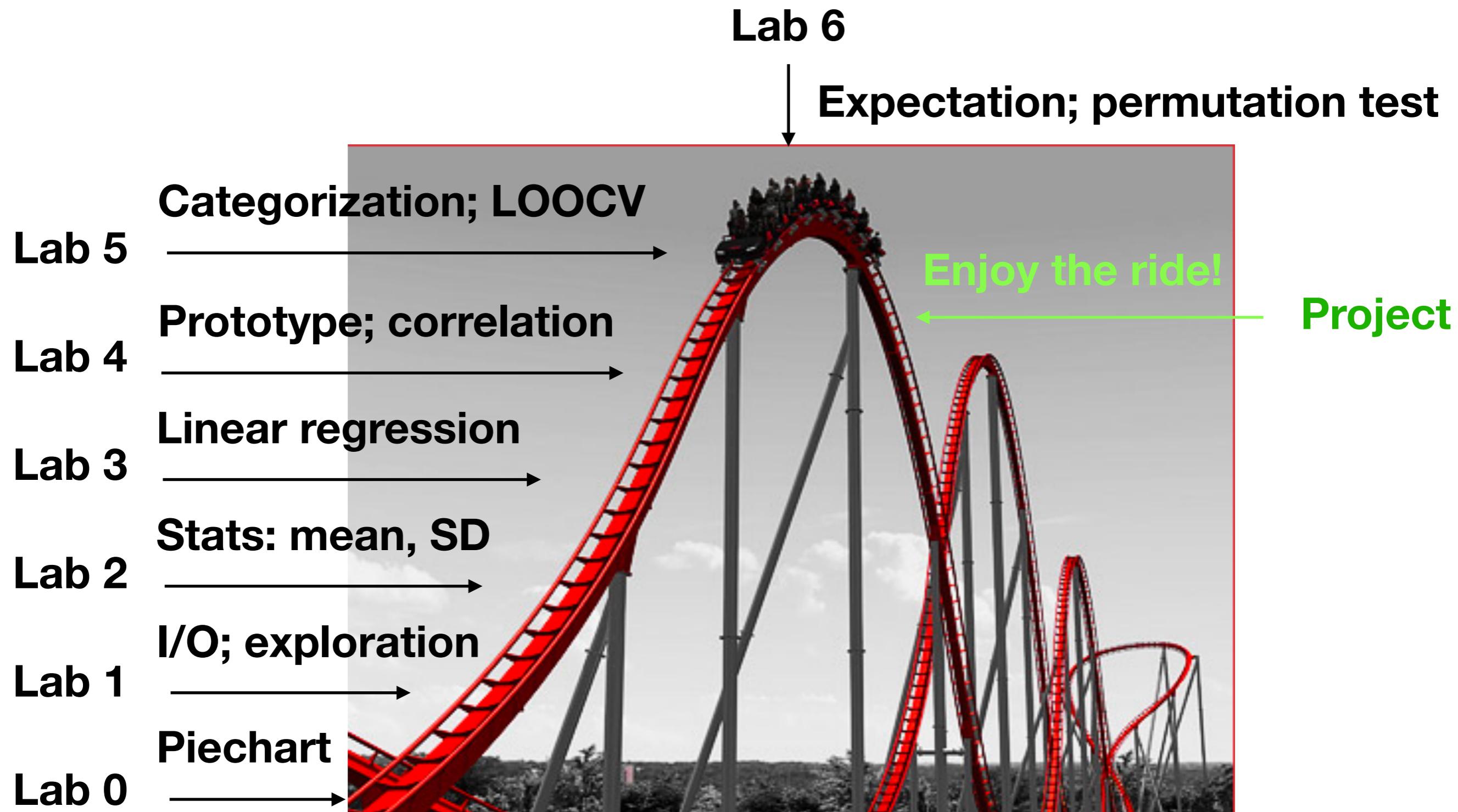
Up to now



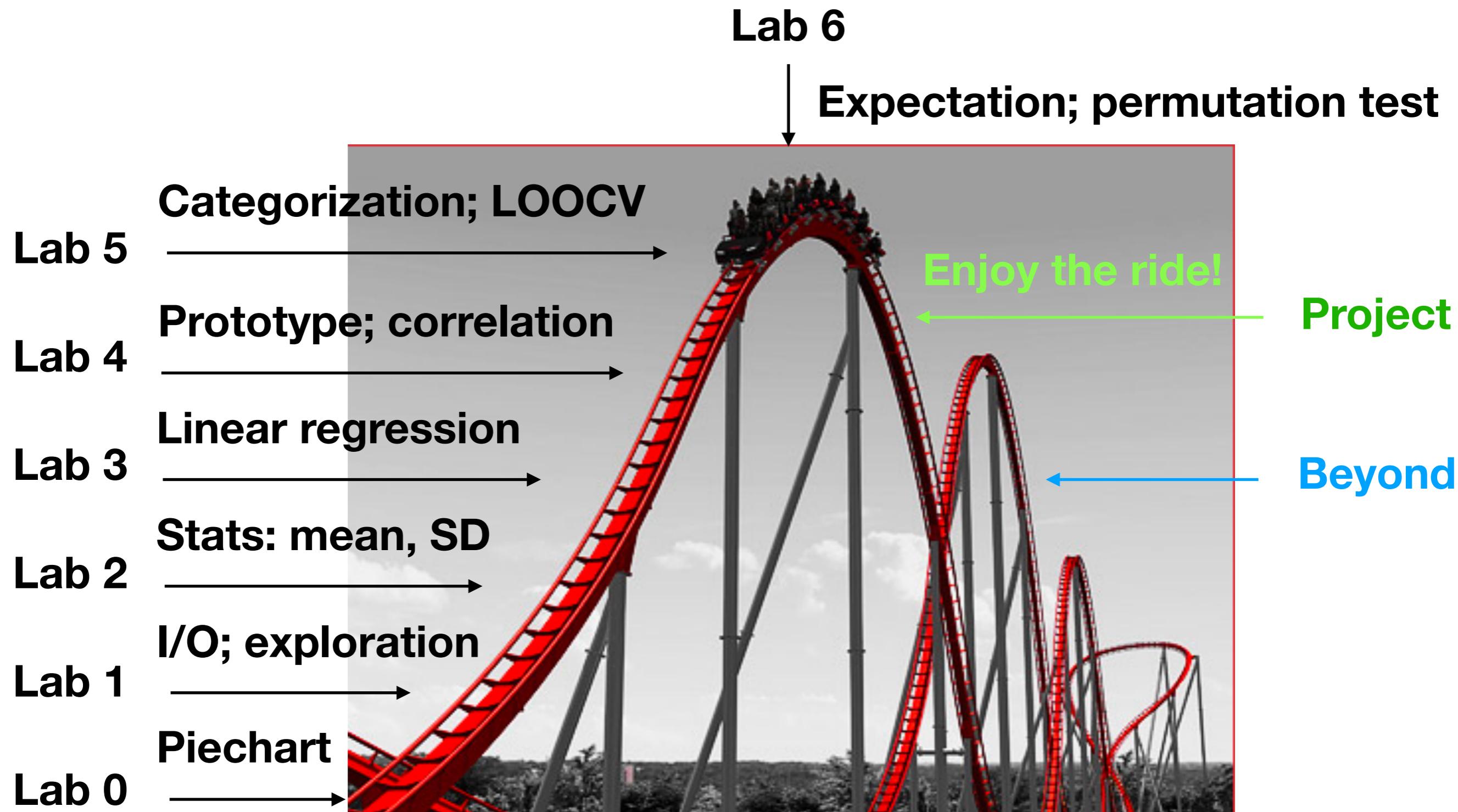
Up to now



Up to now



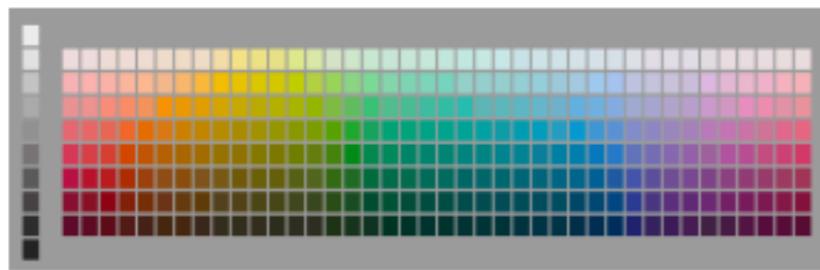
Up to now



Outline

- Cross-linguistic variation and universals
- Course project

Color



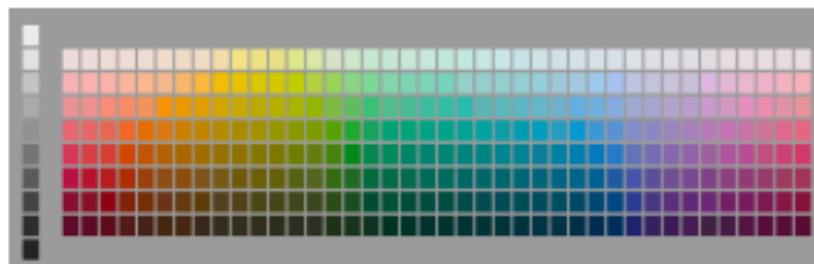
Wobé (Ivory Coast)



Buglere (Panama)



Color



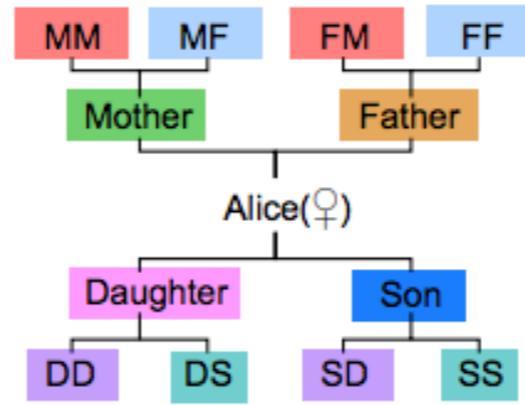
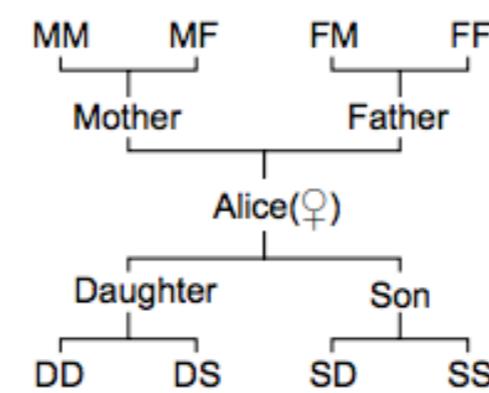
Wobé (Ivory Coast)



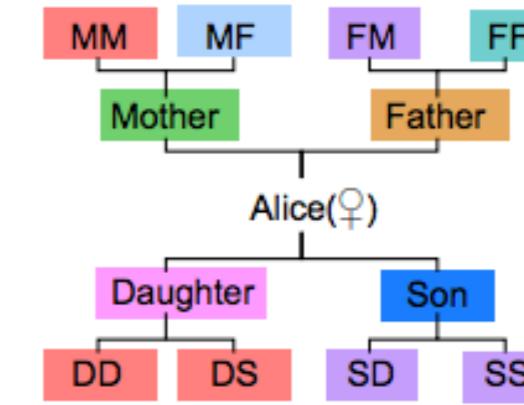
Buglere (Panama)



Kinship

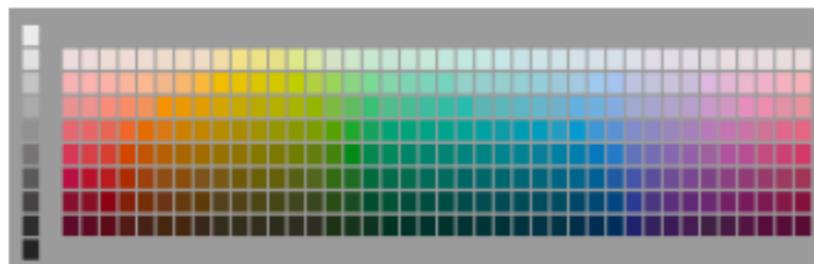


English



Northern Paiute (California)

Color

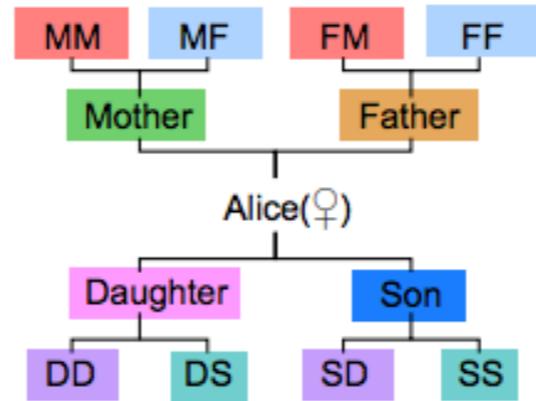
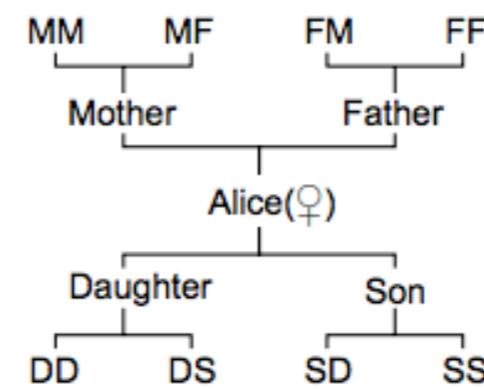


Wobé (Ivory Coast)

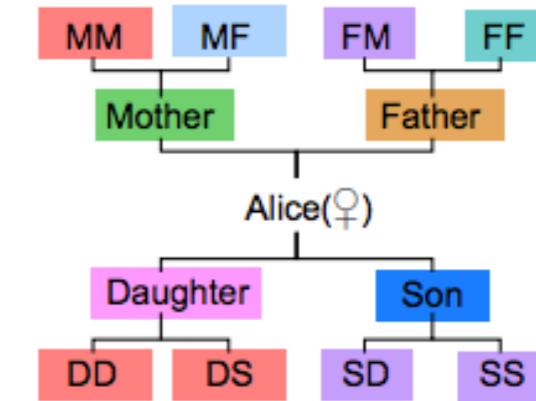


Buglere (Panama)

Kinship



English



Northern Paiute (California)

Space

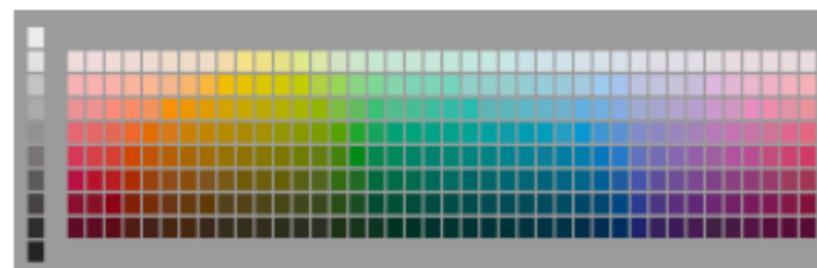


English



Chichewa (Malawi)

Color



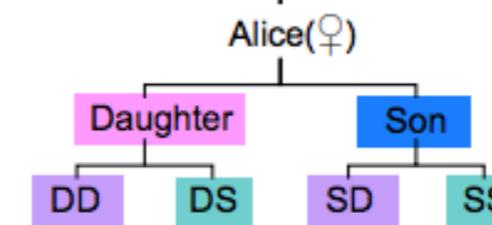
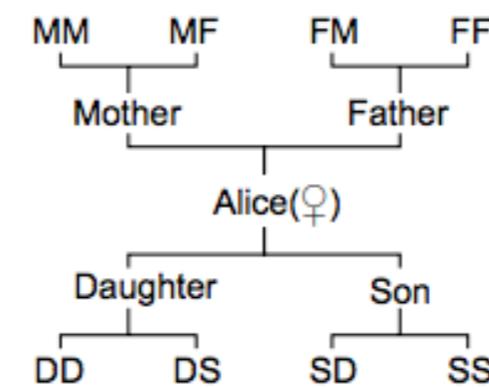
Wobé (Ivory Coast)



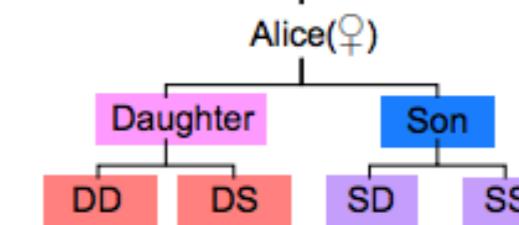
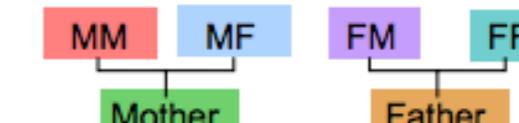
Buglere (Panama)



Kinship



English



Northern Paiute (California)

Space



Number 1 2 3 4 5 6 7 8 9

English



Chichewa (Malawi)



English



Piraha



Some fundamental questions

- Are there “universal” constraints across languages?
 - Linguistic universals

Some fundamental questions

- Are there “universal” constraints across languages?
 - Linguistic universals
 - How do language specificities influence cognition?
 - Linguistic relativity (or the Sapir-Whorf hypothesis)

Some fundamental questions

- Are there “universal” constraints across languages?
 - Linguistic universals
 - How do language specificities influence cognition?
 - Linguistic relativity (or the Sapir-Whorf hypothesis)
 - Why do languages vary the way they do?
 - Linguistic typology
-

Some fundamental questions

- Are there “universal” constraints across languages?
 - Linguistic universals
- How do language specificities influence cognition?
 - Linguistic relativity (or the Sapir-Whorf hypothesis)
- Why do languages vary the way they do?
 - Linguistic typology
-

Color

The world color survey (WCS)



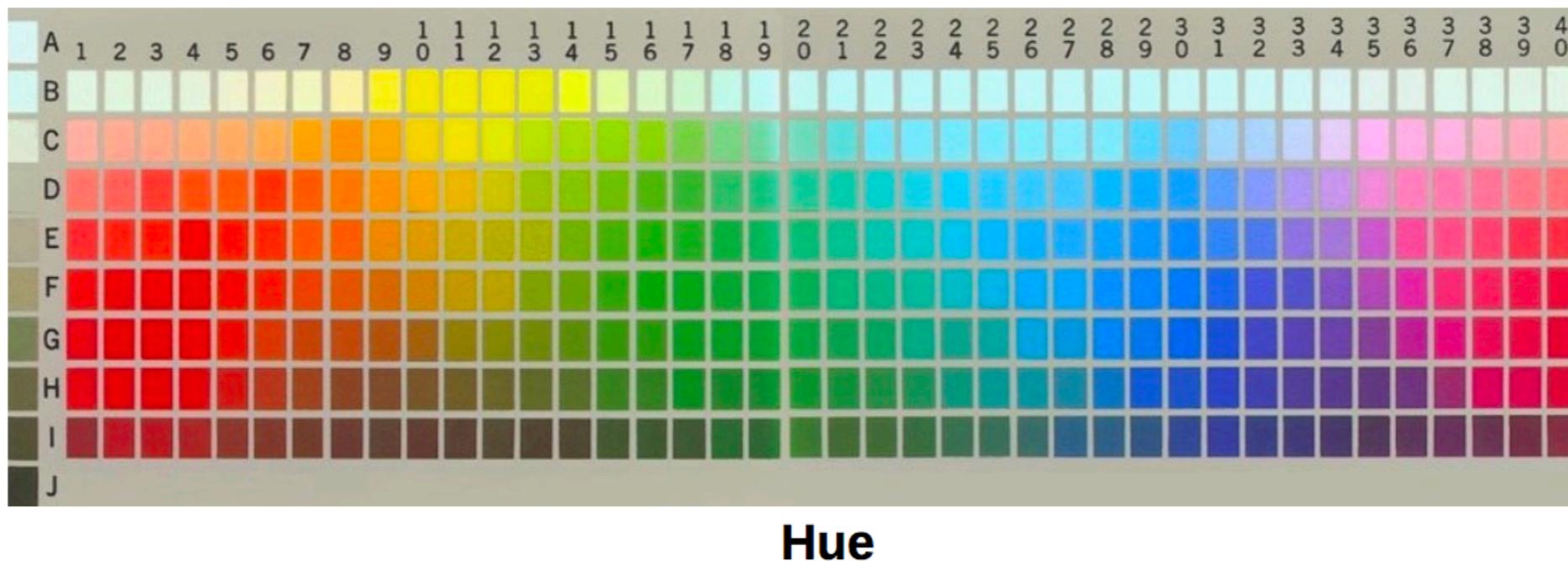
The World Color Survey

The World Color Survey (WCS) was initiated in the late 1970's to test the hypotheses advanced by Berlin and Kay ([1969](#)) regarding

- (1) the existence of universal constraints on cross-language color naming, and
- (2) the existence of a partially fixed evolutionary progression according to which languages gain color terms over time.

330 Munsell color stimuli at maximal chroma used in the World Color Survey.

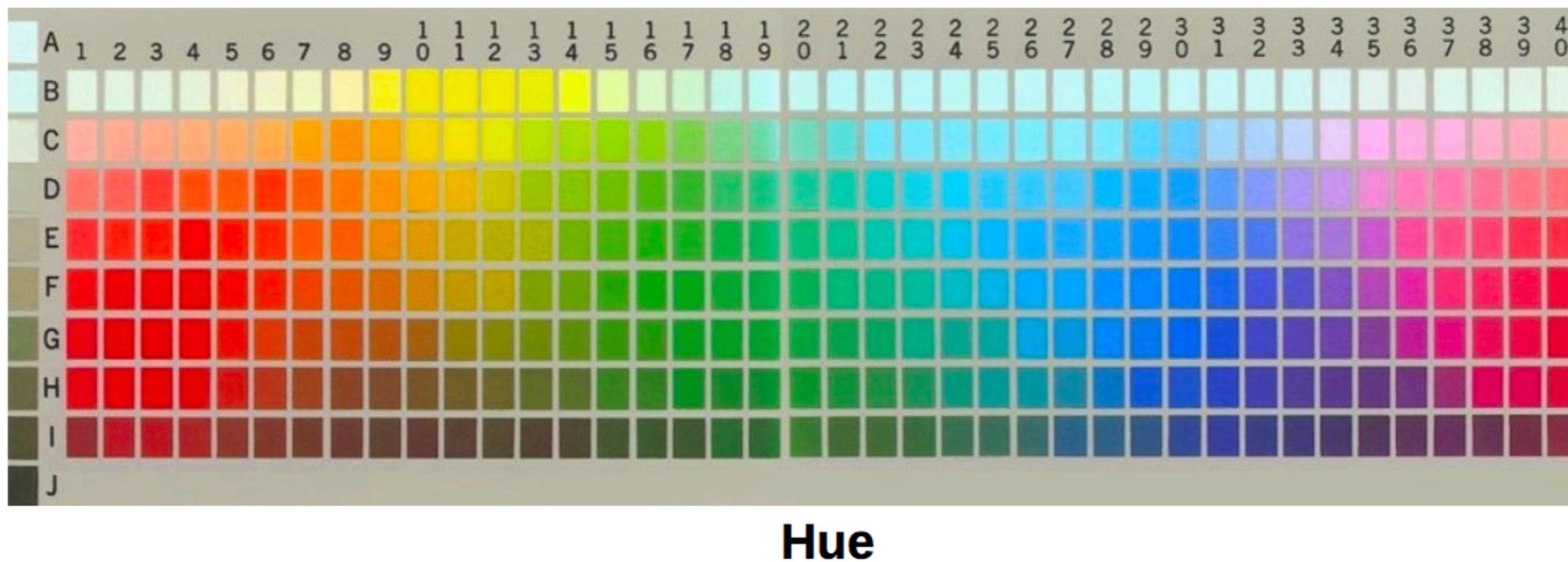
Lightness



Hue

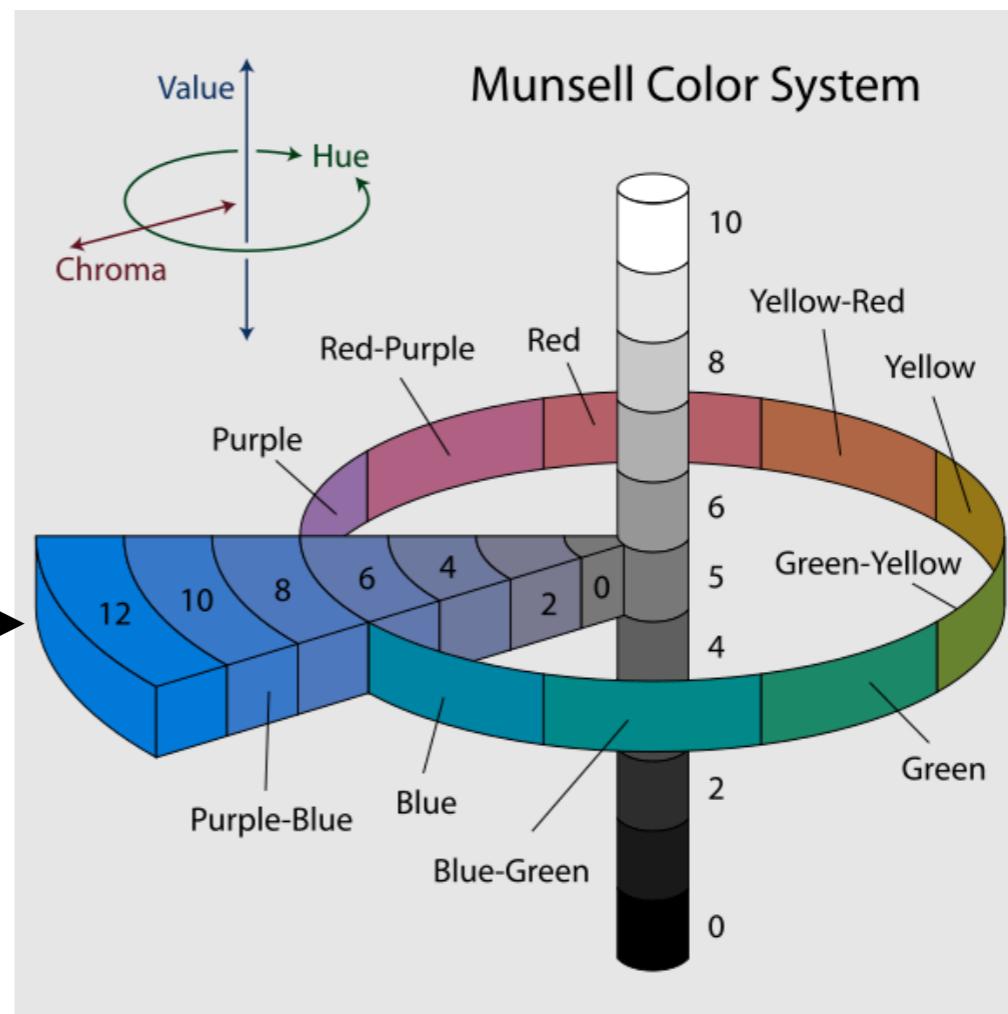
330 Munsell color stimuli at maximal chroma used in the World Color Survey.

Lightness

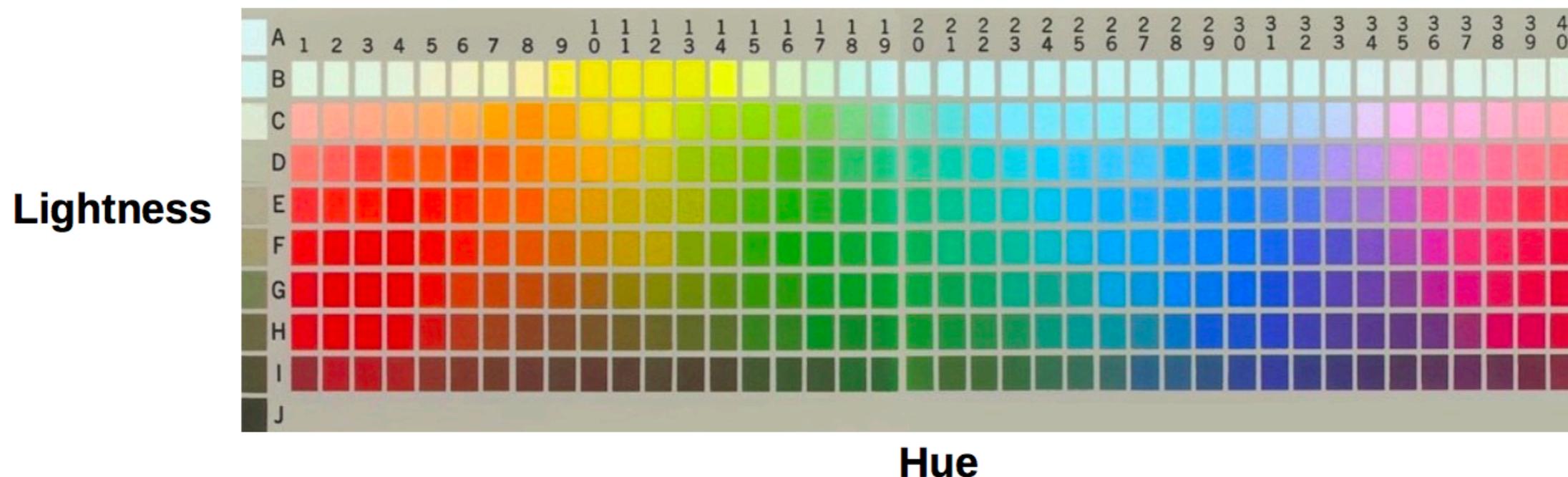


Hue

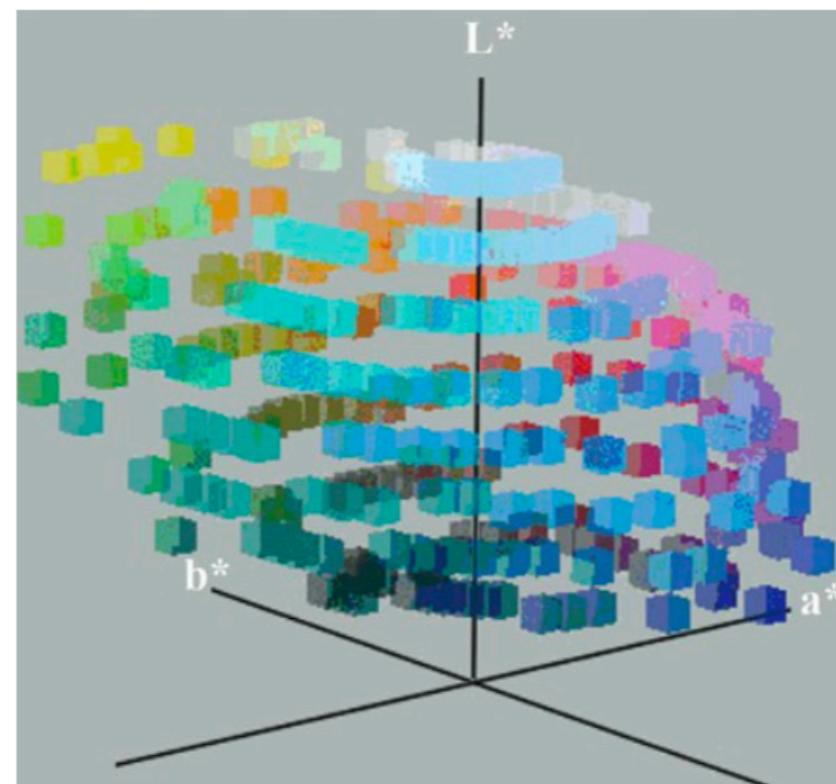
The chroma dimension
is roughly constant in WCS



330 Munsell color stimuli at maximal chroma used in the World Color Survey.



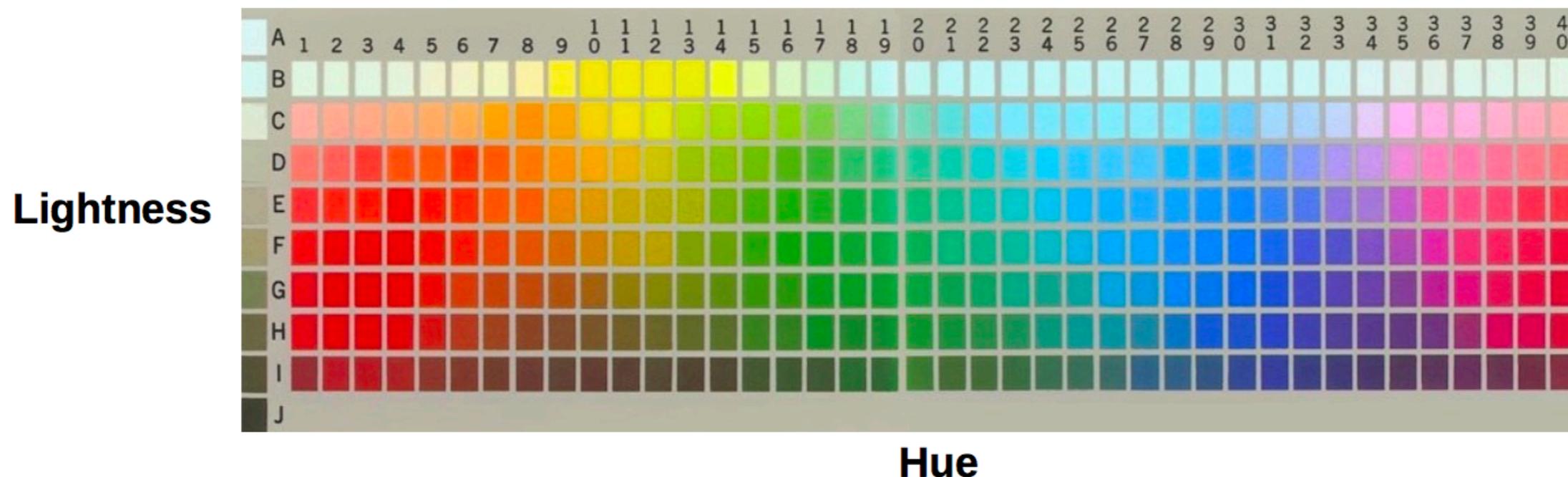
Same stimuli visualized in CIELAB (~perceptual) space.



L: Lightness
AB: a-b color opponency

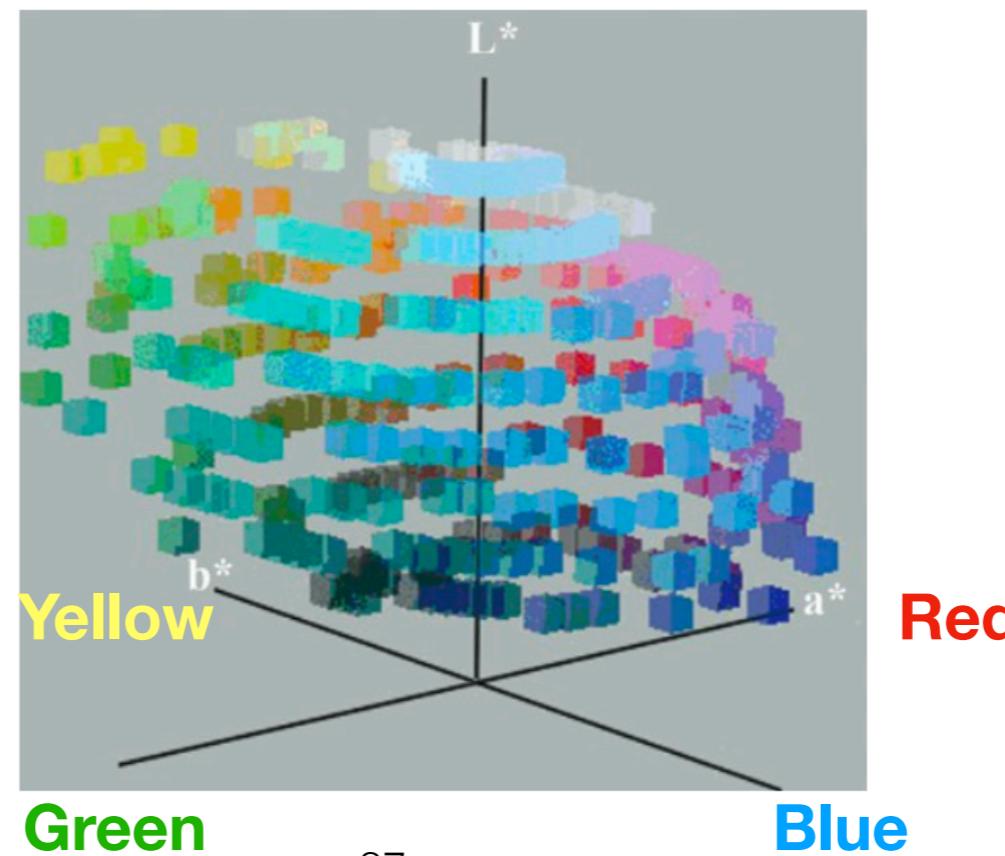
An alternative color space,
also available in WCS.

330 Munsell color stimuli at maximal chroma used in the World Color Survey.



Same stimuli visualized in CIELAB (~perceptual) space.

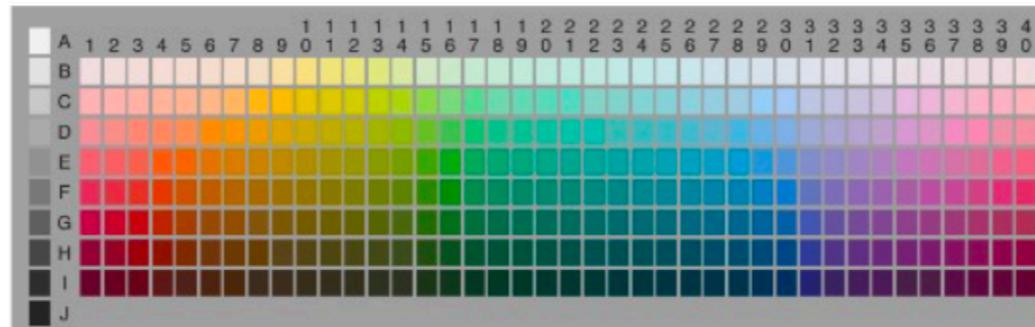
L: Lightness
AB: a-b color opponency



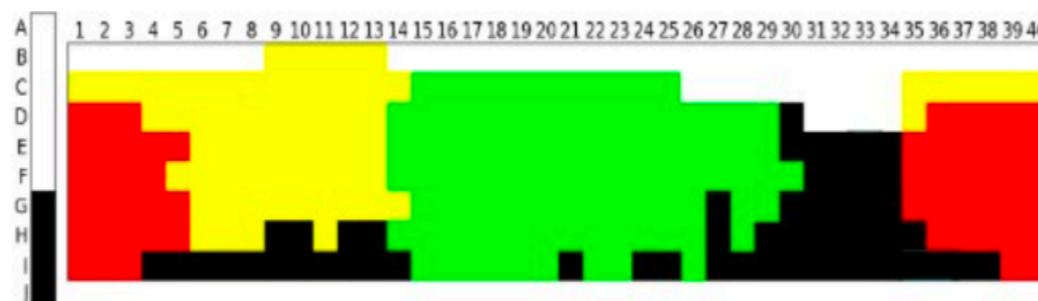
Opponency 1:
Green vs Red

Opponency 2:
Yellow vs Blue

Cross-language variation in color naming



Stimulus array

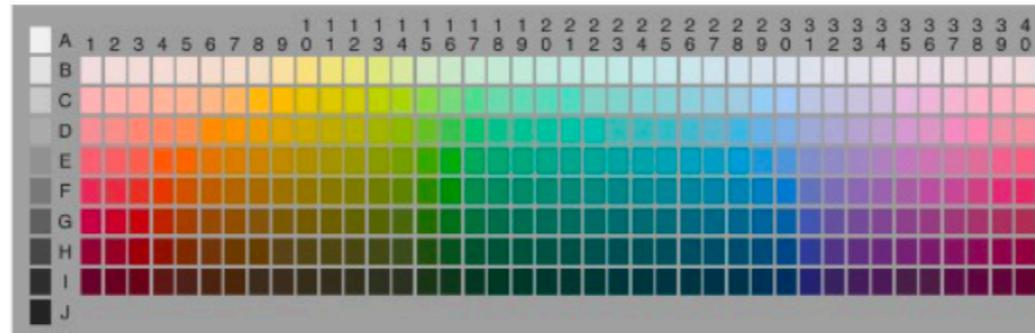


Iduna (Papua New Guinea)

Each color patch represents a “partition” from a modal color term in a given language.

Adapted from Regier, Kemp & Kay (2015)

Cross-language variation in color naming



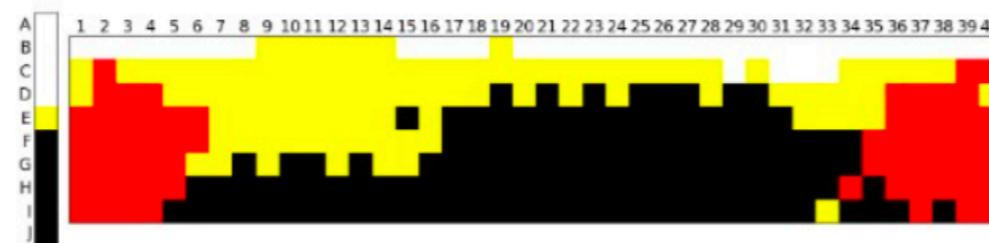
Stimulus array

K = 3



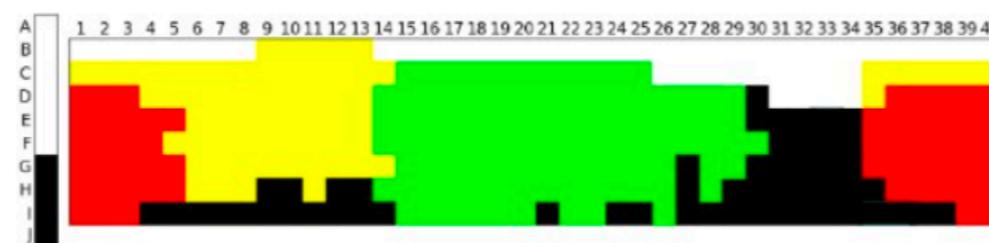
Bantoid (Nigeria/Cameroon)

K = 4



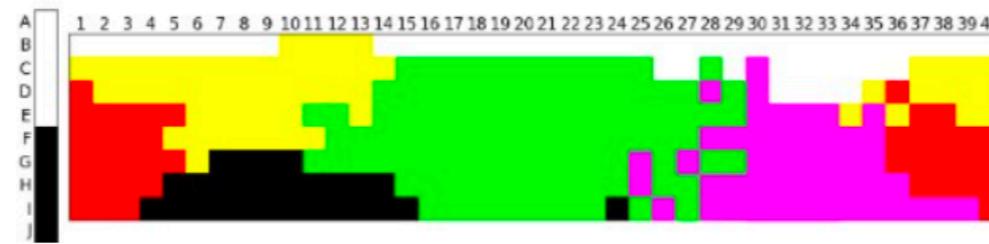
Culina (Peru/Brazil)

K = 5



Iduna (Papua New Guinea)

K = 6



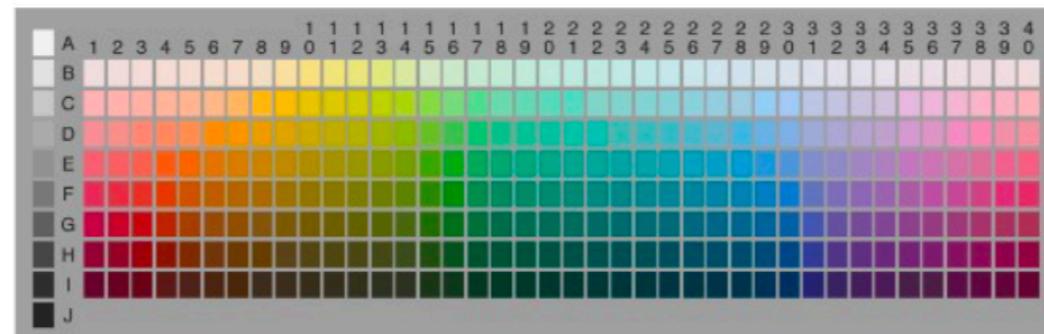
Buglere (Panama)

Adapted from Regier, Kemp & Kay (2015)

A classic research question

- Are there “universal” constraints in color naming across different languages?

Color foci



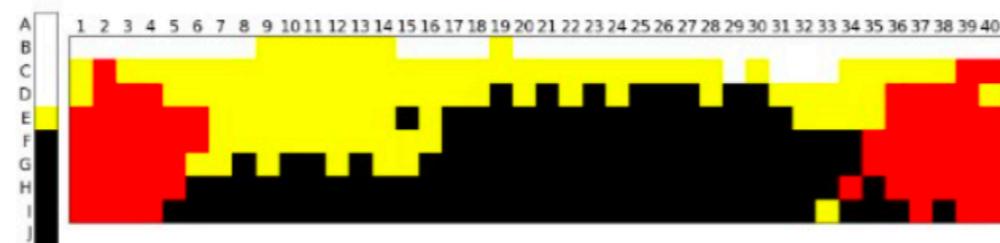
Stimulus array

K = 3



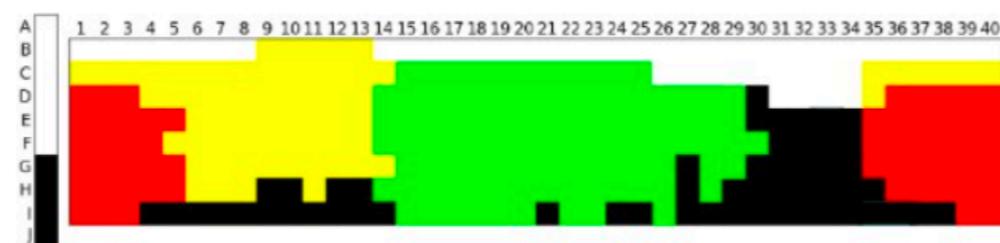
Bantoid (Nigeria/Cameroon)

K = 4



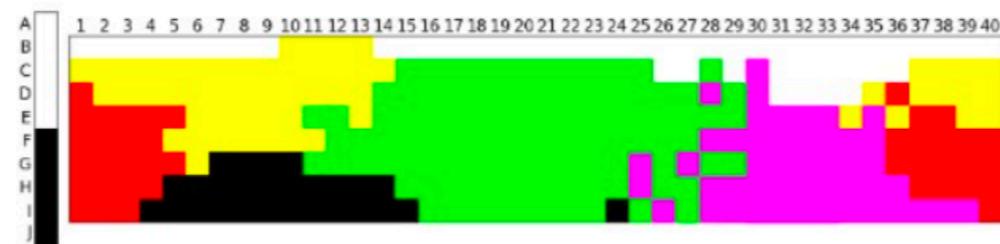
Culina (Peru/Brazil)

K = 5



Iduna (Papua New Guinea)

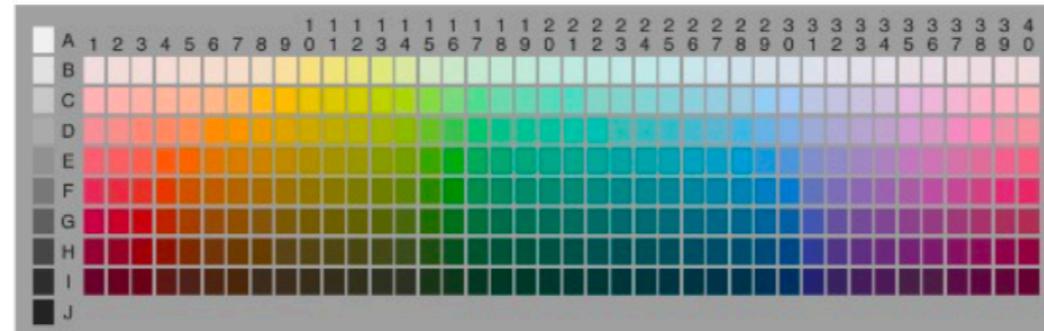
K = 6



Buglere (Panama)

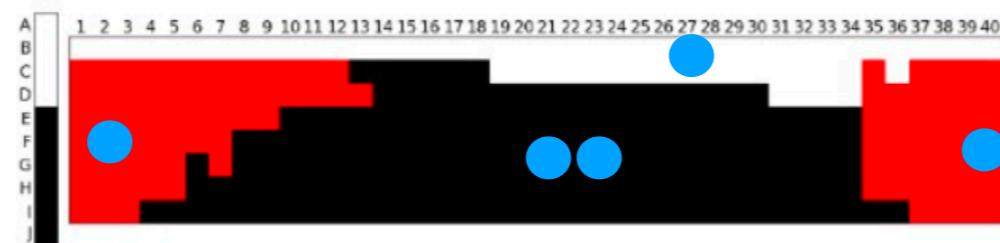
Adapted from Regier, Kemp & Kay (2015)

Color foci



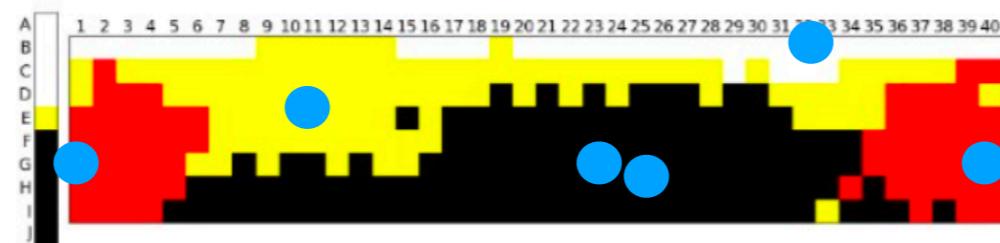
Stimulus array

K = 3



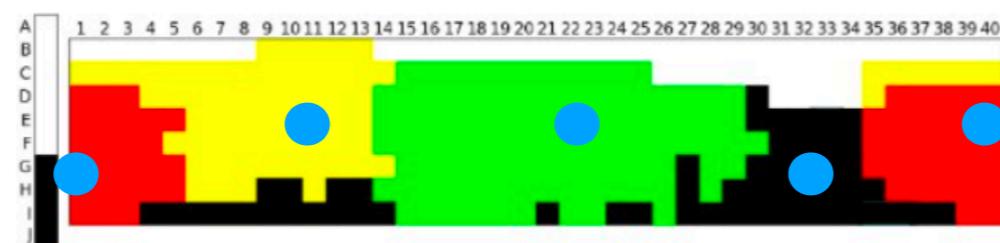
Bantoid (Nigeria/Cameroon)

K = 4



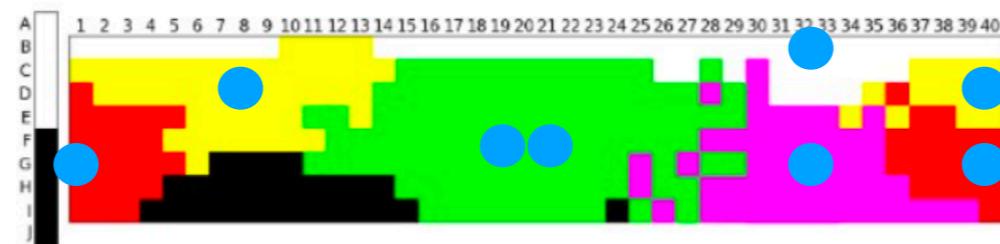
Culina (Peru/Brazil)

K = 5



Iduna (Papua New Guinea)

K = 6

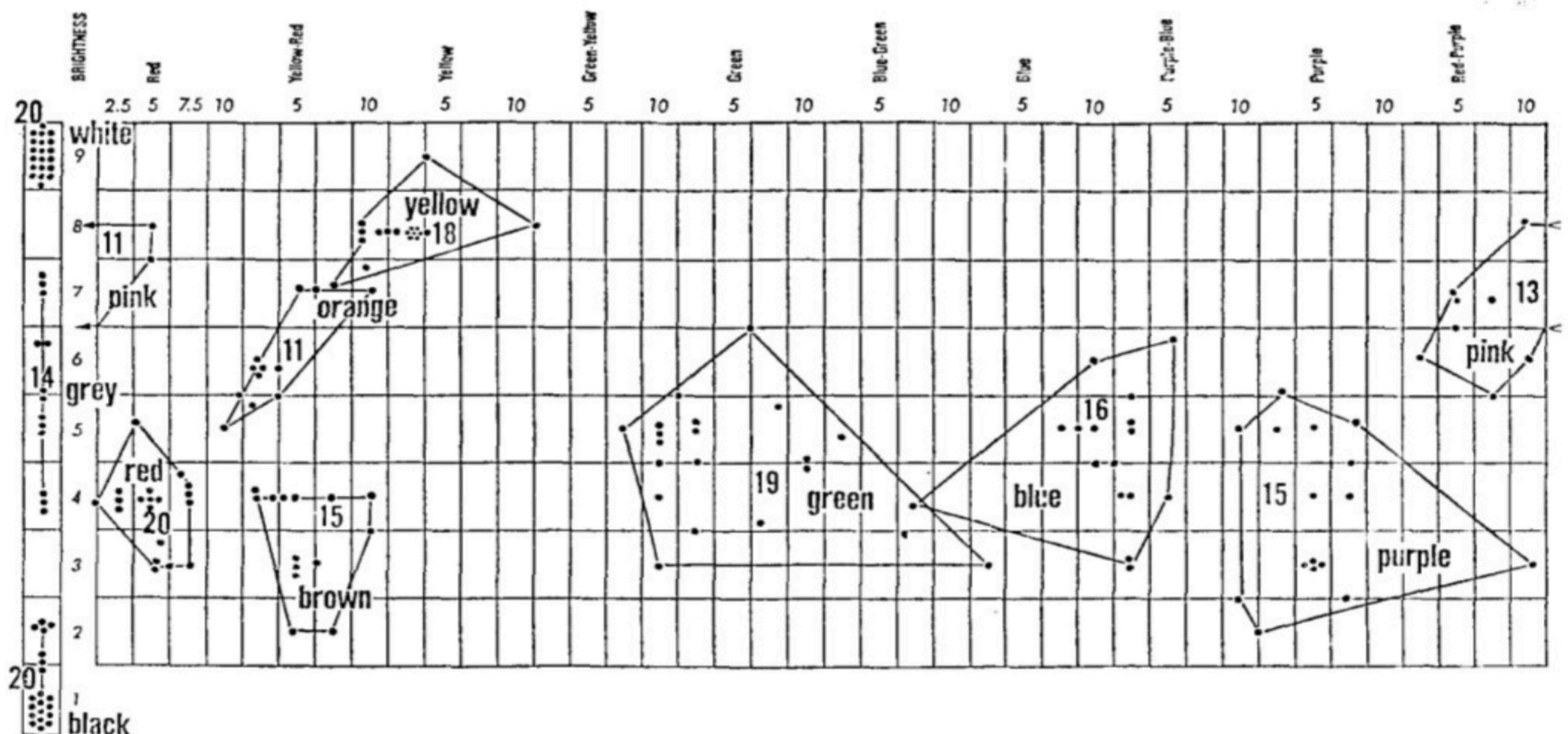


Buglere (Panama)

Adapted from Regier, Kemp & Kay (2015)

Constraint I (Berlin & Kay, 1969)

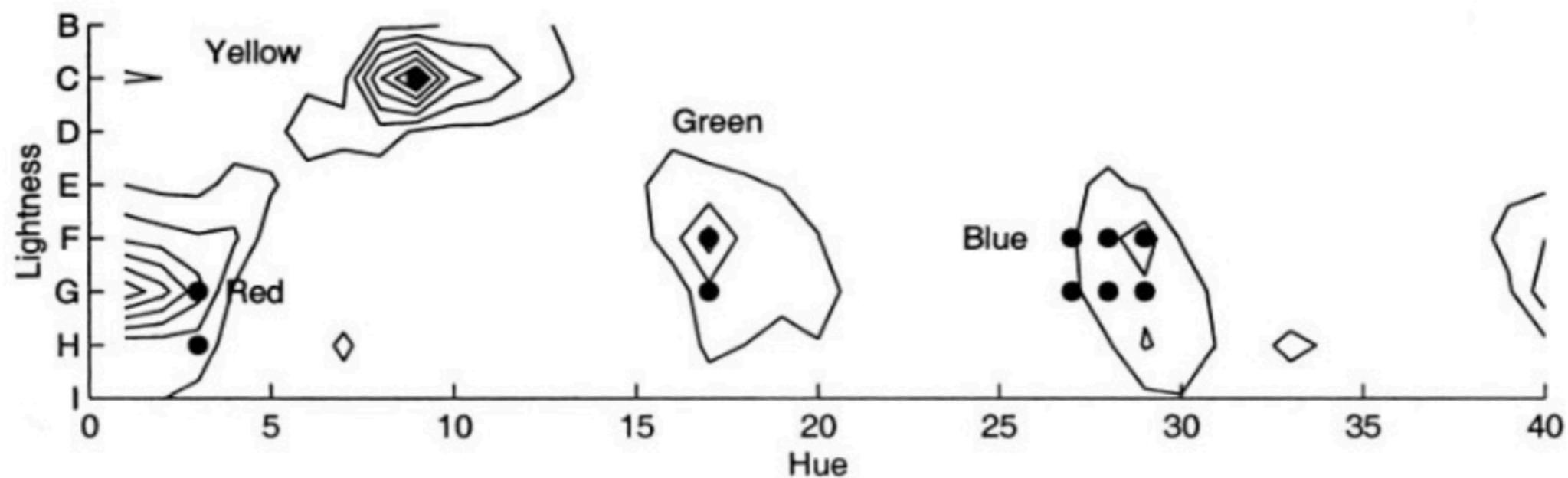
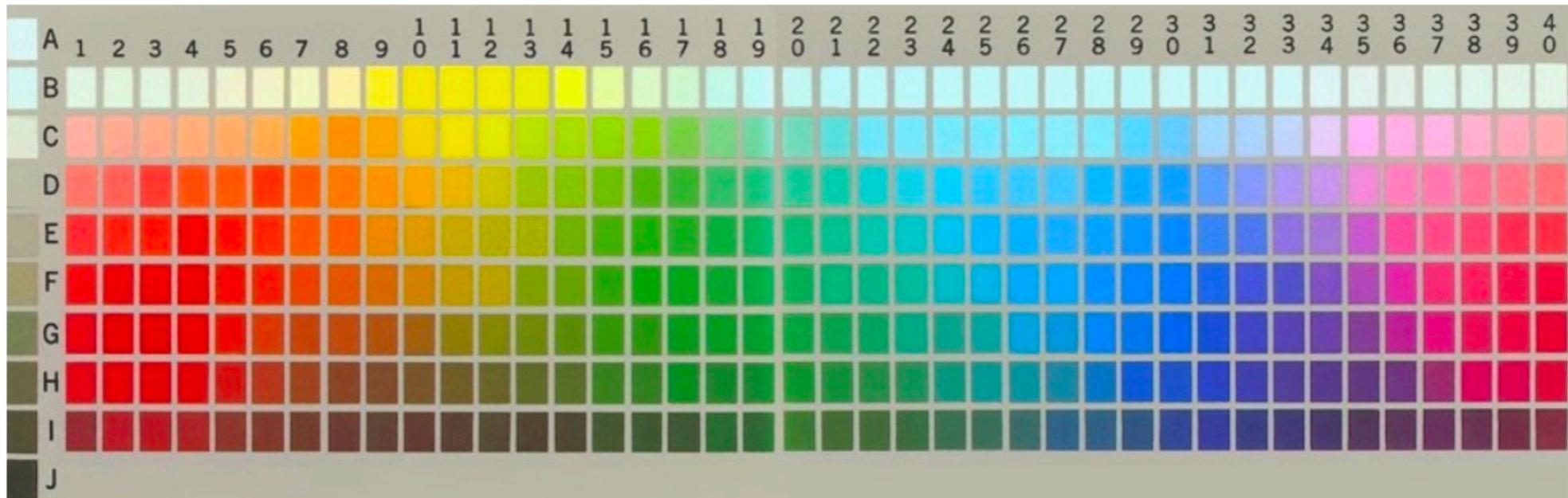
- Color foci (or best examples) are clustered x-lg.



Data from 20 languages (B&K, 1969)

Constraint I

- Further evidence from 110 languages in the World Color Survey.



Data from 110 languages (Regier, Kay & Cook, 2005)

Constraint II (Berlin & Kay, 1969)

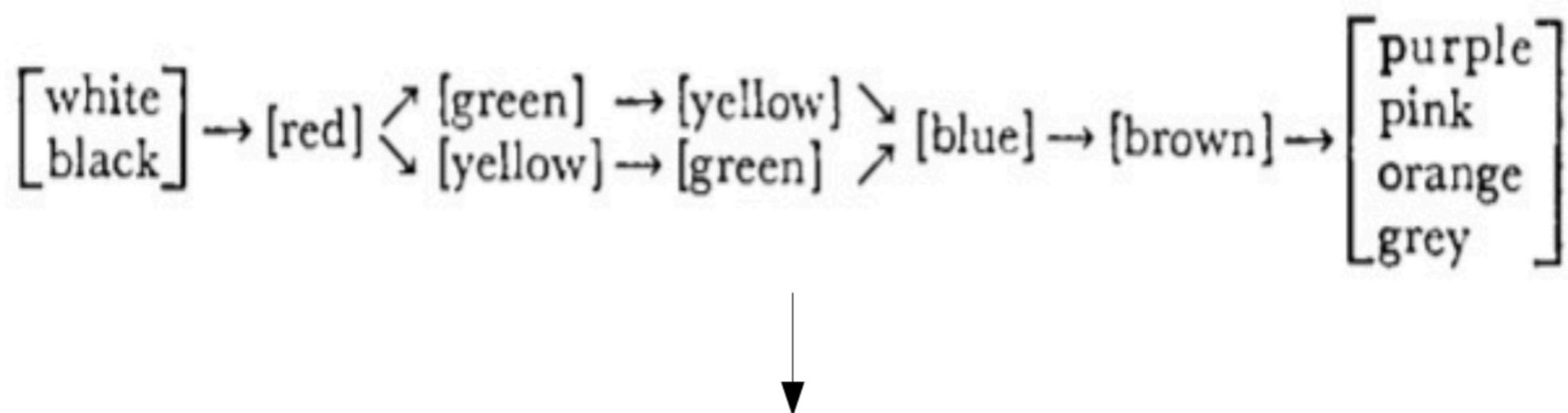
THE TWENTY-TWO ACTUALLY OCCURRING TYPES
OF BASIC COLOR LEXICON

Type	No. of basic color terms	Perceptual categories encoded in the basic color terms										
		white	black	red	green	yellow	blue	brown	pink	purple	orange	grey
1	2	+	+	-	-	-	-	-	-	-	-	-
2	3	+	+	+	+	-	-	-	-	-	-	-
3	4	+	+	+	+	+	-	-	-	-	-	-
4	4	+	+	+	+	+	-	-	-	-	-	-
5	5	+	+	+	+	+	-	-	-	-	-	-
6	6	+	+	+	+	+	-	-	-	-	-	-
7	7	+	+	+	+	+	-	-	-	-	-	-
8	8	+	+	+	+	+	-	-	-	-	-	-
9	8	+	+	+	+	+	-	-	-	-	-	-
10	8	+	+	+	+	+	-	-	-	-	-	-
11	8	+	+	+	+	+	-	-	-	-	-	-
12	9	+	+	+	+	+	-	-	-	-	-	-
13	9	+	+	+	+	+	-	-	-	-	-	-
14	9	+	+	+	+	+	-	-	-	-	-	-
15	9	+	+	+	+	+	-	-	-	-	-	-
16	9	+	+	+	+	+	-	-	-	-	-	-
17	9	+	+	+	+	+	-	-	-	-	-	-
18	10	+	+	+	+	+	-	-	-	-	-	-
19	10	+	+	+	+	+	-	-	-	-	-	-
20	10	+	+	+	+	+	-	-	-	-	-	-
21	10	+	+	+	+	+	-	-	-	-	-	-
22	11	+	+	+	+	+	-	-	-	-	-	-

NOTE: Only these twenty-two out of the logically possible 2,048 combinations of the eleven basic color categories are found.

Constraint II (Berlin & Kay, 1969)

- Regularity in color names implies an evolutionary path.

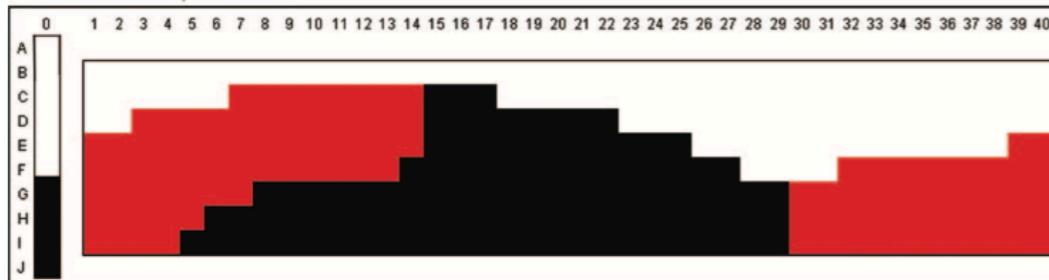


- Black and White (Bk&W): Distinguish black and white.
- Warm and Cool (Wa&C): Distinguish the warm primaries (red and yellow) from the cool primaries (green and blue).
- Red: Distinguish red.

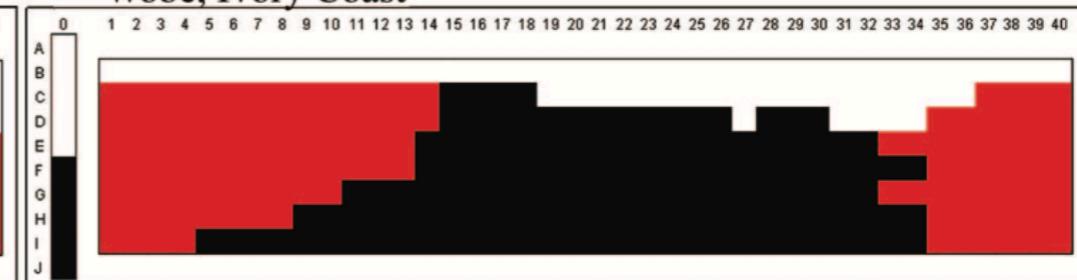
Modification by Kay & Maffi (1999) based on WCS 110 lgs: Still contested.

Near-optimal color naming

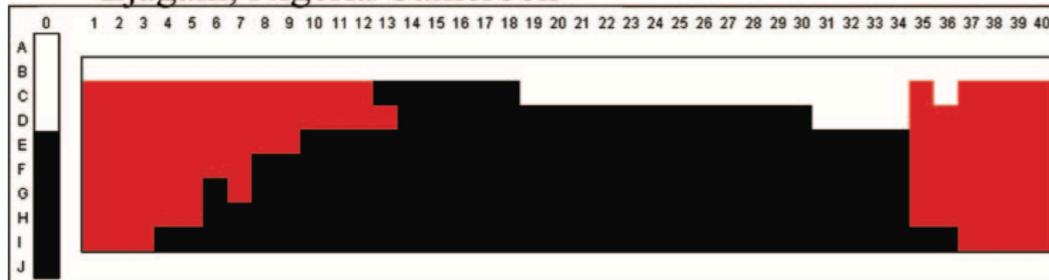
Model, n=3



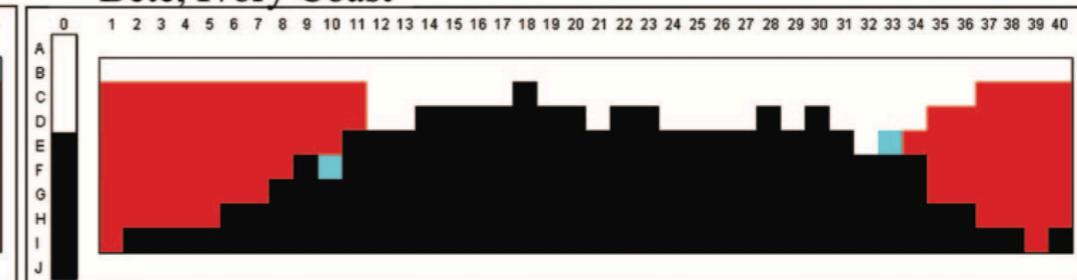
Wobé, Ivory Coast



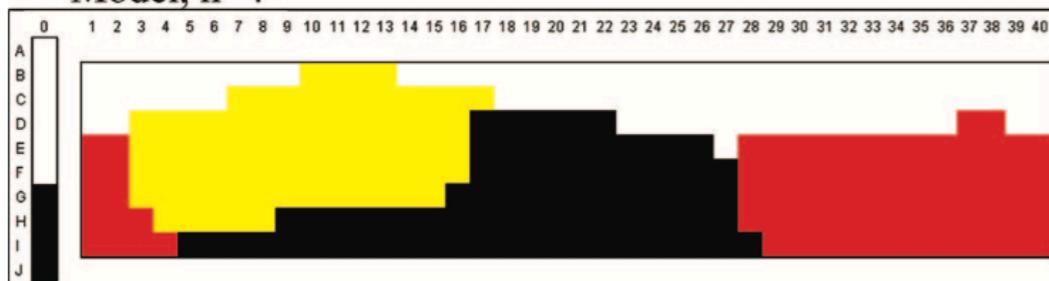
Ejagam, Nigeria/Cameroon



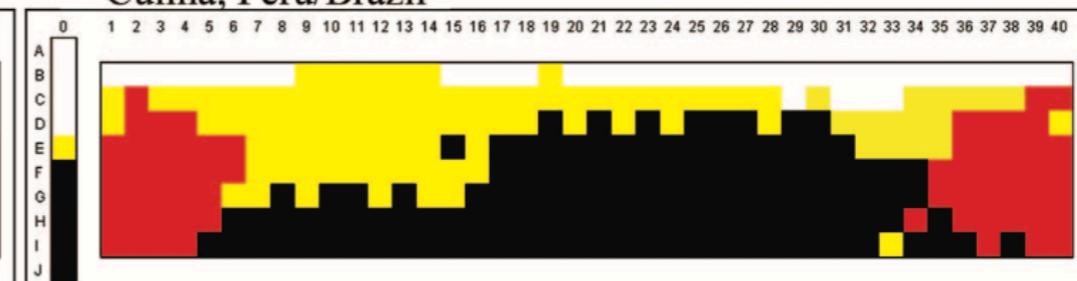
Bété, Ivory Coast



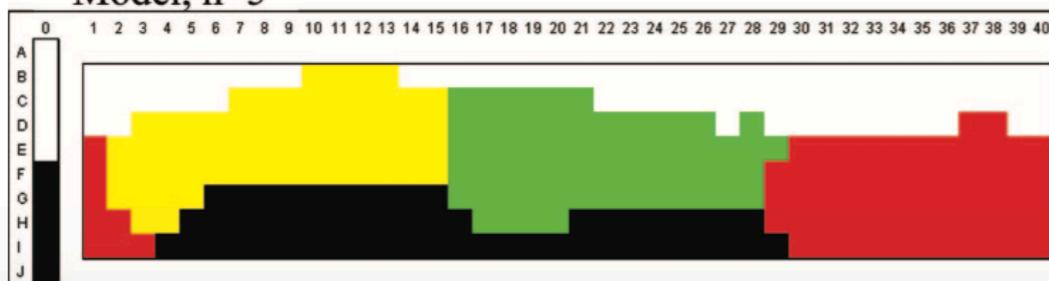
Model, n=4



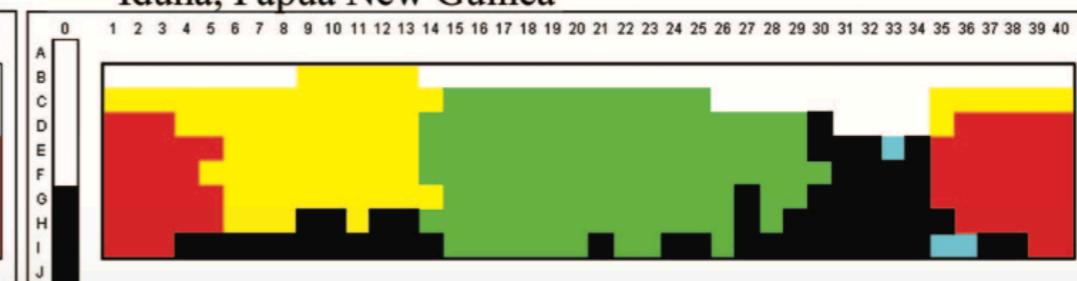
Culina, Peru/Brazil



Model, n=5



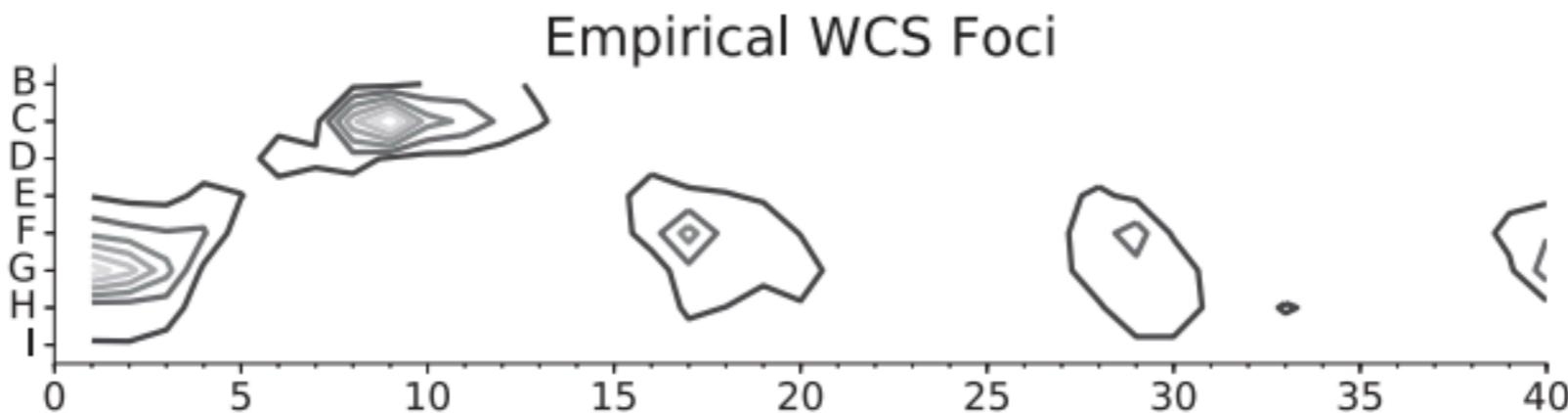
Iduna, Papua New Guinea



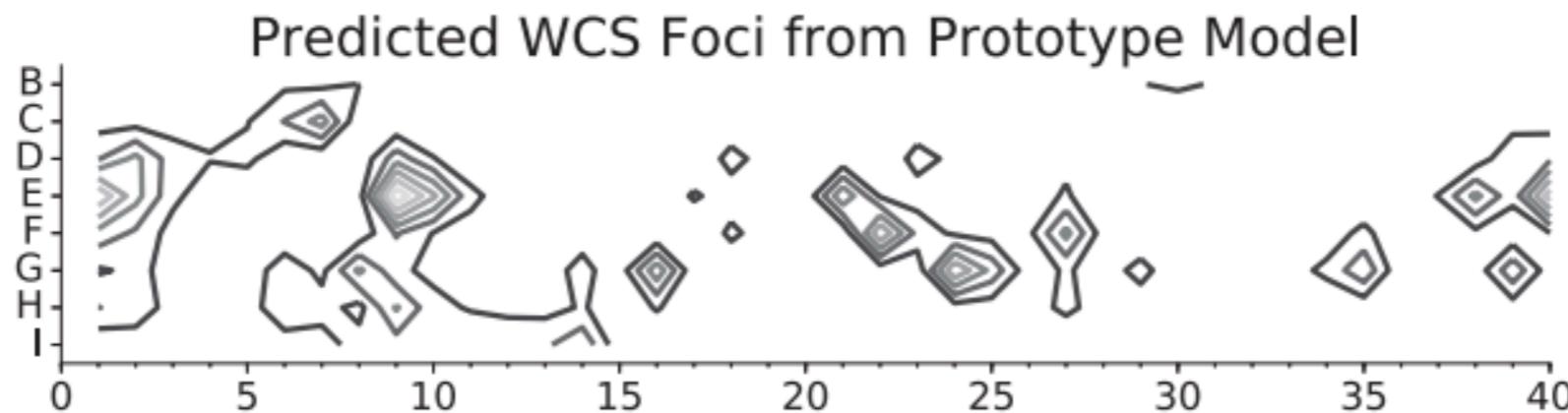
Regier, Kay, & Khetarpal (2007); Kemp, Xu, & Regier (2018)

Models of color foci

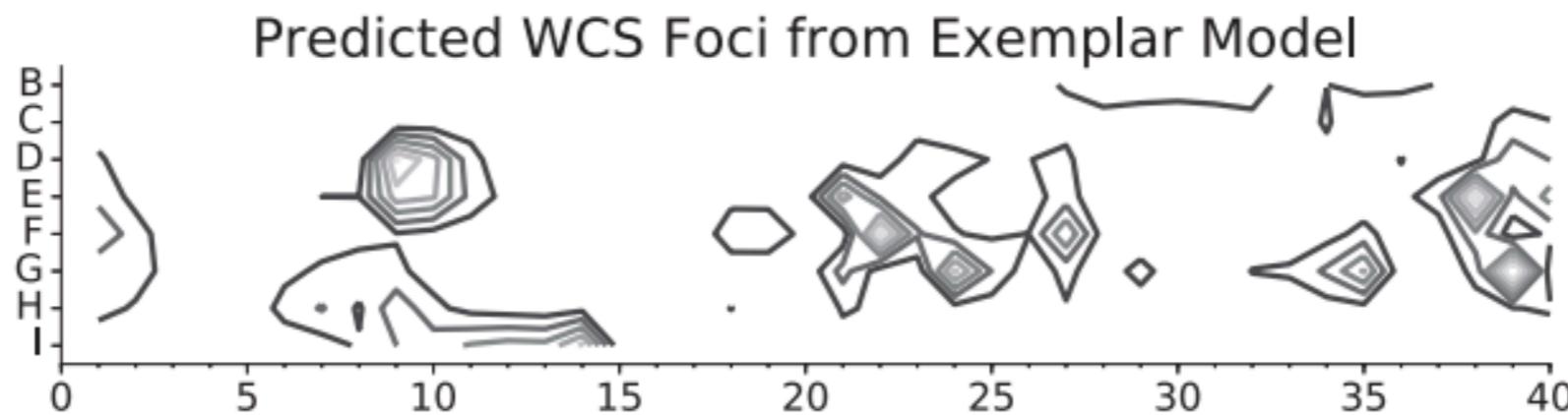
A



D



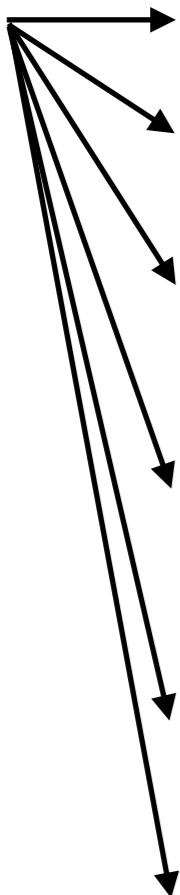
E



Abbott, Griffiths, & Regier (2016)

Reading materials

Project-related (WCS)



Required reading:

- Berlin, B., and Kay, P. (1969). *Basic color terms: Their universality and evolution*. University of California Press.

Optional readings:

- Regier, T., Kay, P., and Cook, R. S. (2005). Focal colors are universal after all. *Proceedings of the National Academy of Sciences*, 102(23), 8386–8391.
- Regier, T., Kay, P., and Khetarpal, N. (2007). Color naming reflects optimal partitions of color space. *Proceedings of the National Academy of Sciences*, 104(4), 1436–1441.
- Abbott, J. T., Griffiths, T. L., and Regier, T. (2016). Focal colors across languages are representative members of colors categories. *Proceedings of the National Academy of Sciences*, 113(40), 11178-11183.
- Gibson, E., Futrell, R., Jara-Ettinger, J., Mahowald, K., Bergen, L., Ratnasingam, S., ... and Conway, B. R. (2017). Color naming across languages reflects color use. *Proceedings of the National Academy of Sciences*, 114(40), 10785–10790.
- Kemp, C., Xu, Y., and Regier, T. (2018). Semantic typology and efficient communication. *Annual Review of Linguistics*, 4, 109–128.
- Zaslavsky, N., Kemp, C., Regier, T., and Tishby, N. (2018). Efficient compression in color naming and its evolution. *Proceedings of the National Academy of Sciences*, 115(31), 7937–7942.

Recommended book:

- Pinker, S. (2007). *The stuff of thought: Language as a window into human nature*. Penguin.

How does color naming influence cognition?



RESEARCH ARTICLE

The Sapir-Whorf Hypothesis and Probabilistic Inference: Evidence from the Domain of Color

Emily Cibelli¹✉, Yang Xu^{2,3}, Joseph L. Austerweil⁴, Thomas L. Griffiths^{3,5}, Terry Regier^{2,3*}

1 Department of Linguistics, Northwestern University, Evanston, IL 60208, United States of America,

2 Department of Linguistics, University of California, Berkeley, CA 94720, United States of America,

3 Cognitive Science Program, University of California, Berkeley, CA 94720, United States of America,

4 Department of Psychology, University of Wisconsin, Madison, WI 53706, United States of America,

5 Department of Psychology, University of California, Berkeley, CA 94720, United States of America



5-minute break

Plan for course-project weeks

- Next week: Judgment and decision making + Project
- Next week x 2: Gender + Project
- Next week x 3: Morality + Project
- Final-class week: Data Blitz! (project report due on that Friday)

Course project

- Goal: Collaborative, hypothesis-driven exploratory data analysis
- Time frame: ~1 month (including in-class project time)
- Data: 1) world color survey, or 2) a choice of your own *
- Starting point: Background research on previous work on color naming
- Assessments (ALL instructions on syllabus via Quercus):
 - Proposal (5%)
 - Data blitz (10%)
 - Final report (15%)

Project instructions

In this project, you will work with your collaborator to formulate a cognitive hypothesis and test it against an extensive public data set. The default data set you will be working with is [The World Color Survey](#) or WCS ([manual](#)), which includes color naming and behavioural data from 110 non-industrialized languages. You have the option of working individually and/or with an alternative data set of your choice that is similar in complexity to WCS. If you choose to do so, you will need to obtain permission from the instructor *prior to the project proposal due date*.

For the WCS data analysis, you may work with the demo Jupyter Notebook provided. For analysis with data sets sourced elsewhere, you are responsible to develop your own Jupyter Notebook. All analyses should be performed in Python Jupyter Notebook. It is sufficient to turn in a **single copy** of the proposal and the final report between you and your collaborator (specify your names in the front page). Late submission will receive a deduction of 1 point per delayed hour. Submit your proposal and report in PDF format through Quercus. No other laboratory work will be assigned during the project period.

- All available on the syllabus page

Project instructions

Stage I. Initial proposal (due 1 week since announcement) [5pts]

Submit a PDF proposal (about 2 pages) that describes a concrete plan for your project. Your proposal should include 1) a description of a specific and testable cognitive hypothesis [1pt] 2) proposed methodologies, models, and analyses [2pts] 3) expected results such as any figures and/or tables you will generate from the proposed analyses [1.5pt] 4) division of labour between you and your collaborator, if you work in pairs [.5pt]. Your proposal needs to be feasible given the time frame of the project. Invalid or infeasible proposals will potentially delay your progress; they will be returned and revised.

Stage II. Data blitz (final day of class) [10pts]

You and your collaborator are required to present in the data blitz. It is recommended that you split the presentation into two even halves since you will be assessed individually, not as a group. Slides in PDF format should be submitted as PDF through the online submission system at least 1 day prior to the presentation. Failure to conform to the required PDF format will receive 2 points deduction; delayed submission of your slides will incur the same penalty as any late assignment. During the data blitz session, presentation order will be randomized. Failure to attend the data blitz will receive 0 credit. Presentation should be kept under 5-6 minutes and will be timed strictly.

- All available on the syllabus page

Project instructions

Stage III. Final report (due the Monday after data blitz) [15pts]

Submit a PDF document of no more than 10 pages (excluding Appendix for Python code) that includes:

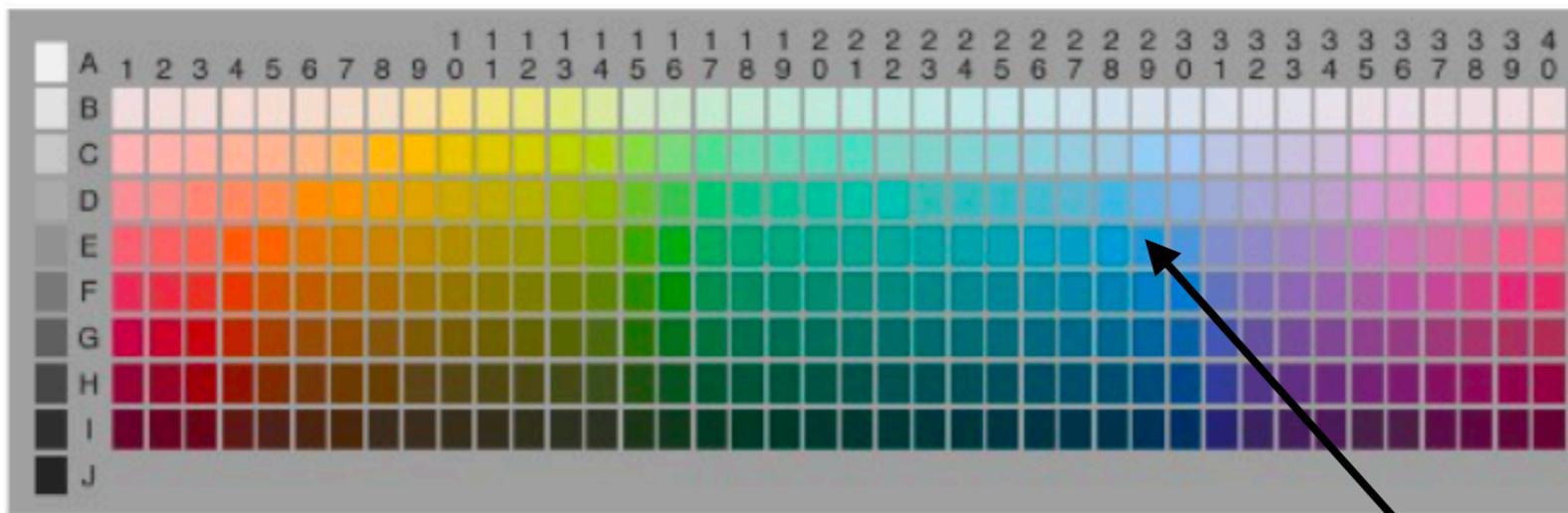
- 1) *Abstract*: In 1 paragraph, summarize your hypothesis, methods, findings, and conclusions. [1pt]
- 2) *Introduction*: Provide background on the project. Justify and articulate your hypothesis. Specify at least 1 alternative hypothesis. [3pts]
- 2) *Methods*: Present the methodologies for your analyses in a way that is reproducible. Provide the full set of Python code in Appendix that allows reproduction of your reported findings. [4pts]
- 3) *Results*: Summarize the main results from your analyses by making use of figures, tables, and statistics. Explain how these findings support or refute your hypothesis. [3pts]
- 4) *Conclusion*: Conclude in 1-2 paragraphs; discuss possible extensions. [2pt]
- 5) *References*: List at least 5 full references with in-text citations. [2pts]

- All available on the syllabus page

World color survey

- URL: <http://www1.icsi.berkeley.edu/wcs/>
- READ the manual and browse the data before you begin
- Notebook primer provided on Syllabus page
- Data archive: <http://www1.icsi.berkeley.edu/wcs/data.html>
 - Primary WCS data (included in Notebook primer folder)
 - Berlin & Kay (1969) data (optional)
 - WCS mapping tables (auxiliary)

“WCS_data_core”: “chip.txt”



Stimulus array

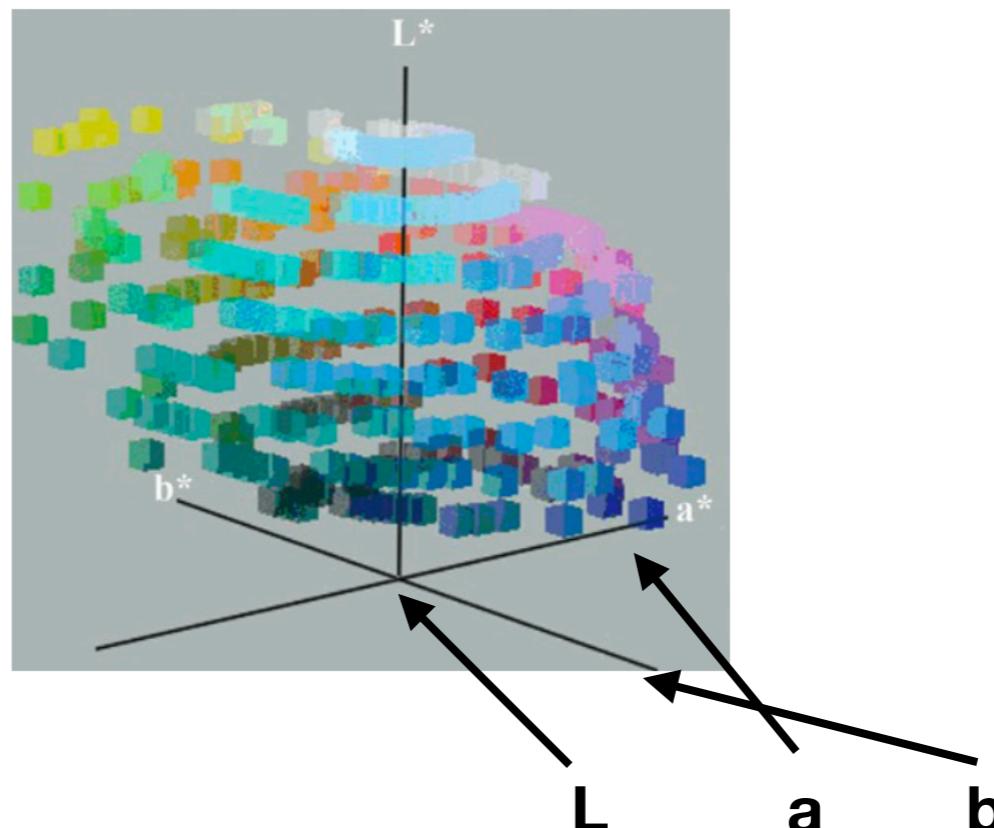
Chip Index

1	E	29	E29
2	C	23	C23
3	F	4	F4
4	I	36	I36
5	C	20	C20
6	C	6	C6
7	E	15	E15
8	H	40	H40
9	G	6	G6
10	I	30	I30

“cnum-vhcm-lab-new.txt”

Same stimuli visualized in CIELAB (~perceptual) space.

L: Lightness
AB: a-b color opponency



Chip
Index

141	A	0	0	10.00RP	9.5	96.00	-.06	.06
274	B	0	0	10.00RP	9	91.08	-.05	.06
129	B	1	2	2.50R	9	91.08	5.53	2.22
230	B	2	2	5.00R	9	91.08	5.51	3.28
302	B	3	2	7.50R	9	91.08	5.54	4.46
324	B	4	2	10.00R	9	91.08	5.43	5.64
27	B	5	2	2.50YR	9	91.08	5.21	7.67
290	B	6	2	5.00YR	9	91.08	4.30	10.08

“spkr-lsas.txt”

Language Index (1-110)	Speaker Index	Age	Gender
1	1	90	M
1	2	26	M
1	3	38	M
1	4	35	M
1	5	80	M
1	6	48	M
1	7	26	M
1	8	39	M
1	9	47	F
1	10	49	M
1	11	40	F
1	12	45	M
1	13	50	M
1	14	30	M
1	15	21	M
1	16	60	F
1	17	32	M
1	18	67	M
1	19	15	M
1	20	42	M
1	21	40	M
1	22	47	M
1	23	23	F
1	24	45	F
1	25	30	F
2	1	20	F
2	2	40	F
2	3	45	F

“term.txt”

Language Index (1-110)	Speaker Index	Chip Index	Term Gloss (abr)
1	1	1	LB
1	1	2	LB
1	1	3	LE
1	1	4	WK
1	1	5	LF
1	1	6	LE
1	1	7	F
1	1	8	LE
1	1	9	LE
1	1	10	LB
1	1	11	LB
1	1	12	F
1	1	13	LB
1	1	14	LB
1	1	15	LF
1	1	16	LF
1	1	17	LE

“foci-exp.txt”



Bantoid (Nigeria/Cameroon)

Language Index (1-110)	Speaker Index	Term Index	Term Gloss	Foci Location
1	1	1	LF	A0
1	1	2	WK	D9
1	1	2	WK	D10
1	1	2	WK	D11
1	1	2	WK	D12
1	1	3	F	D25
1	1	4	LB	J0
1	1	5	G	F17
1	1	6	LE	F1
1	1	6	LE	F2
1	1	6	LE	F3
1	1	6	LE	G1
1	1	6	LE	G2
1	1	6	LE	G3
1	2	1	LF	A0
1	2	2	S	C1
1	2	2	S	C2
1	2	2	S	C3
1	2	2	S	C4
1	2	2	S	C5
1	2	3	WK	E1

Project possibilities

- Replication of universal color naming analyses, e.g. B&K (1969), or Regier, Kay, & Cook (2005)

Project possibilities

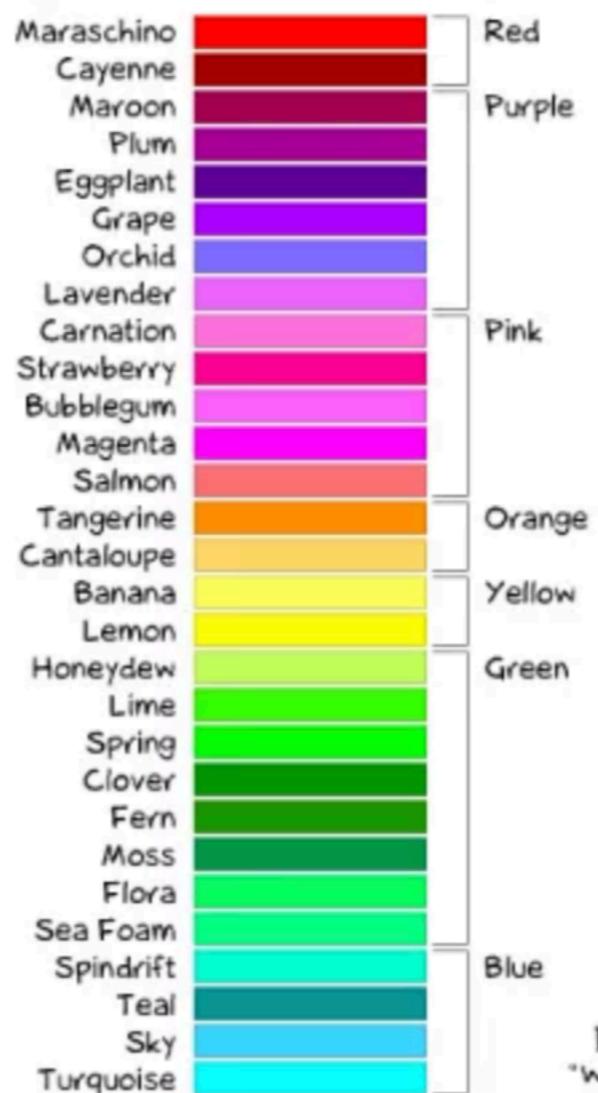
- Replication of universal color naming analyses, e.g. B&K (1969), or Regier, Kay, & Cook (2005)
- Question 1: Why does color naming vary across languages, i.e., what “caused” cross-language variation?

Project possibilities

- Replication of universal color naming analyses, e.g. B&K (1969), or Regier, Kay, & Cook (2005)
- Question 1: Why does color naming vary across languages, i.e., what “caused” cross-language variation?
- Question 2: Is there a gender difference in color naming?

Gender differences in color?

Color names if
you're a girl...



Color names if
you're a guy...

Doghouse Diaries
"We take no as an answer."

Gender and color?

*Actual color names
if you're a girl ...*



*Actual color names
if you're a guy ...*

Basically, women were slightly more liberal with the modifiers

Project possibilities

- Replication of universal color naming analyses, e.g. B&K (1969), or Regier, Kay, & Cook (2005)
- Question 1: Why does color naming vary across languages, i.e., what “caused” cross-language variation?
- Question 2: Is there a gender difference in color naming?
- Question 3: Is there an age difference in color naming?

Project possibilities

- Replication of universal color naming analyses, e.g. B&K (1969), or Regier, Kay, & Cook (2005)
- Question 1: Why does color naming vary across languages, i.e., what “caused” cross-language variation?
- Question 2: Is there a gender difference in color naming?
- Question 3: Is there an age difference in color naming?
- Question 4: How are color categories represented in the mind, e.g., are names of the color chips predictable from prototype, exemplar, or other categorization models? What is the “best” way of naming color?

Project possibilities

- Replication of universal color naming analyses, e.g. B&K (1969), or Regier, Kay, & Cook (2005)
- Question 1: Why does color naming vary across languages, i.e., what “caused” cross-language variation?
- Question 2: Is there a gender difference in color naming?
- Question 3: Is there an age difference in color naming?
- Question 4: How are color categories represented in the mind, e.g., are names of the color chips predictable from prototype, exemplar, or other categorization models? What is the “best” way of naming color?
- Find and brainstorm with your project collaborator!