

The Feldberg Family Foundation



Carnegie Mellon



mozilla

SIGIR
Special Interest Group
on Information Retrieval

MASCO

YAHOO!

***The Eighth
Annual***

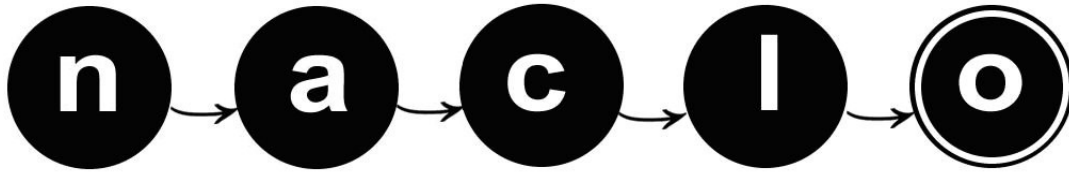
**North American
Computational
Linguistics
Olympiad**

2014

www.naclo.cs.cmu.edu

**Invitational
Round**

March 13, 2014



Welcome to the eighth annual North American Computational Linguistics Olympiad! You are among the few, the brave, and the brilliant, to participate in this unique event. In order to be completely fair to all participants across North America, we need you to read, understand, and follow these rules completely.

Rules

1. The contest is five hours long and includes nine problems, labeled I to Q, divided into two parts—one before lunch and one after.
 2. Follow the facilitators' instructions carefully.
 3. If you want clarification on any of the problems, talk to a facilitator. The facilitator will consult with the jury before answering.
 4. You may not discuss the problems with anyone except as described in items 3 & 12.
 5. Each problem is worth a specified number of points, with a total of 100 points. Make sure to fill out all the answer boxes properly. You are expected to include explanations for most problems in this round.
 6. We will grade only work in this booklet. All your answers should be in the spaces provided in this booklet. **DO NOT WRITE ON THE BACK OF THE PAGES.**
 7. Write your name and registration number on each page:
Here is an example: Jessica Sawyer #850
 8. The top 100 participants (approximately) across the continent in the open round will be invited to the second round.
 9. Each problem has been thoroughly checked by linguists and computer scientists as well as students like you for clarity, accuracy, and solvability. Some problems are more difficult than others, but all can be solved using ordinary reasoning and some basic analytic skills. You don't need to know anything about linguistics or about these languages in order to solve them.
 10. If we have done our job well, very few people will solve all these problems completely in the time allotted. So, don't be discouraged if you don't finish everything.
 11. If you have any comments, suggestions or complaints about the competition, we ask you to remember these for the web-based evaluation. We will send you an e-mail shortly after the competition is finished with instructions on how to fill it out.
 12. **DO NOT DISCUSS THE PROBLEMS UNTIL THEY HAVE BEEN POSTED ONLINE! THIS MAY BE SEVERAL WEEKS AFTER THE END OF THE CONTEST.**
- Oh, and have fun!

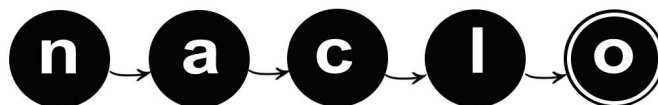
NACLO 2014 Organizers

Program Committee:

Susan Barry, Manchester Metropolitan University
Aleka Blackwell, Middle Tennessee State University
Jordan Boyd-Graber, University of Maryland
Bozhidar Bozhanov, Ontotext
Alan Chang, Princeton University
John DeNero, Google and University of California, Berkeley
Jason Eisner, Johns Hopkins University
Dominique Estival, University of Western Sydney
Matt Gardner, Carnegie Mellon University
Anatole Gershman, Carnegie Mellon University
Linus Hamilton, Massachusetts Institute of Technology
Adam Hesterberg, Massachusetts Institute of Technology
Dick Hudson, University College London
Alex Iriza, Princeton University
Rowan Jacobs, University of Chicago
Wesley Jones, University of Chicago
Mary Laughren, University of Queensland
Lori Levin, Carnegie Mellon University
Patrick Littell, University of British Columbia (co-chair)
Tom McCoy, Yale University
Rachel McEnroe, University of Chicago
David Mortensen, University of Pittsburgh
Babette Newsome, Aquinas College
David Palfreyman, Zayed University
James Pustejovsky, Brandeis University
Dragomir Radev, University of Michigan (co-chair)
Verna Rieschild, Macquarie University
Catherine Sheard, University of Oxford
Ben Sklaroff, University of California, Berkeley
Harold Somers, All Ireland Linguistics Olympiad
Chelsea Voss, Massachusetts Institute of Technology

Problem Credits:

Problem I: Catherine Sheard
Problem J: Tom McCoy
Problem K: David Mortensen
Problem L: Jordan Boyd-Graber
Problem M: David Palfreyman
Problem N: Adam Hesterberg
Problem O: Jonathan Kummerfeld, Aleka Blackwell, and Patrick Littell
Problem P: Jonathan Kummerfeld, Aleka Blackwell, and Patrick Littell
Problem Q: Mary Laughren



NACLO 2014 Organizers (cont'd)

Organizing Committee:

Mary Jo Bensasi, Carnegie Mellon University
Aleka Blackwell, Middle Tennessee State University
Janis Chang, University of Western Ontario
Josh Falk, University of Chicago
Eugene Fink, Carnegie Mellon University
Matt Gardner, Carnegie Mellon University
Adam Hesterberg, Massachusetts Institute of Technology
Alex Iriza, Princeton University
Ann Irvine, Johns Hopkins University
Wesley Jones, University of Chicago
Aaron Klein, Harvard University
Andrew Lamont, Eastern Michigan University
Lori Levin, Carnegie Mellon University (chair)
Jeffrey Ling, Harvard University
Patrick Littell, University of British Columbia
Tom McCoy, Yale University
Rachel McEnroe, University of Chicago
Graham Morehead, University of Maine
David Mortensen, University of Pittsburgh
David Penco, University of British Columbia
James Pustejovsky, Brandeis University
Dragomir Radev, University of Michigan
Julia Workman, University of Montana
Yilu Zhou, Fordham University

Website and Registration:

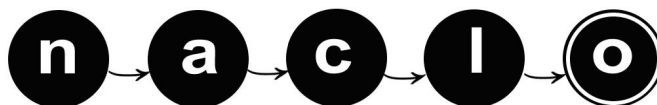
Graham Morehead, University of Maine

US Team Coaches:

Dragomir Radev, University of Michigan (head coach)
Lori Levin, Carnegie Mellon University (coach)

Canadian Coordinator and Team Coach:

Patrick Littell, University of British Columbia



NACLO 2014 Organizers (cont'd)

Contest Site Coordinators:

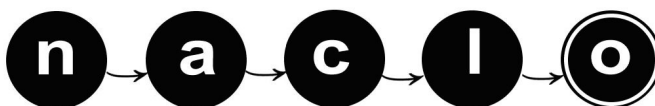
USA

Bemidji State University: Porter Coggins
Brandeis University: James Pustejovsky
Brigham Young University: Deryle Lonsdale
Carnegie Mellon University: Lori Levin, David Mortensen
Central Connecticut State University: Seunghun Lee, Matthew Ciscel, Leyla Zidani-Eroglu
College of William and Mary: Ann Reed
Columbia University: Kathy McKeown, Amy Cooper
Cornell University: Abby Cohn, Sam Tilsen
Dartmouth College: Sravana Reddy
Georgetown University: Daniel Simonson
Johns Hopkins University: Mark Dredze
Massachusetts Institute of Technology: Adam Hesterberg, Chelsea Voss
Middle Tennessee State University: Aleka Blackwell
Minnesota State University, Mankato: Rebecca Bates, Dean Kelley
Northeastern Illinois University: Judith Kaplan-Weinger, Kenneth Konopka
Ohio State University: Micha Elsner, Julie McGory, Michael White
Princeton University: Alan Chang, Christiane Fellbaum, Alex Iriza, Mark Teng
San Jose State University: Hahn Koo, Roula Svorou
Stanford University: Sarah Yamamoto
Stony Brook University: Yejin Choi, Kristen La Magna, Lori Repetti
Union College: Kristina Striegnitz, Nick Webb
University of Alabama at Birmingham: Tamar Solorio
University of Colorado at Boulder: Silva Chang
University of Illinois at Urbana-Champaign: Julia Hockenmaier, Ryan Musa
University of Maine: George Markowsky, Graham Morehead
University of Memphis: Vasile Rus
University of Michigan: Steve Abney, Sally Thomason
University of North Carolina, Charlotte: Wlodek Zadrozny
University of North Texas: Rodney Nielsen, Genevieve Murphy
University of Pennsylvania: Cheryl Hickey
University of Rochester: Mary Swift
University of Southern California: David Chiang
University of Texas: Stephen Wechsler
University of Texas at Dallas: Vincent Ng
University of Washington: Jim Hoard, Luke Zettlemoyer
University of Wisconsin, Madison: Steve Lacy, T.R. Fitz-Gibbon
University of Wisconsin, Milwaukee and Marquette University: Steven Hartman Keiser, Jonas Wittke, Angela Sorby, Hanyong Park, Gabriella Pinter, Joyce Tang Boyland
Western Michigan University: John Kapenga
Western Washington University: Kristin Denham
Yale University: Raffaella Zanuttini, Bob Frank, Aidan Kaplan, Tom McCoy

CANADA

Dalhousie University: Magdalena Jankowska, Vlado Keselj, Armin Sajadi
McGill University: Michael Wagner
Simon Fraser University: Maite Taboada, John Alderete, Cliff Burgess
University of Alberta: Sally Rice
University of British Columbia: Jozina Vander Kloek, David Penco
University of Lethbridge: Yllias Chali
University of Ottawa: Diana Inkpen
University of Toronto: Pen Long, Jordan Ho
University of Western Ontario: Janis Chang

High school sites: Dragomir Radev



NACLO 2014 Organizers (cont'd)

Student Assistants:

Sean Bethard, Brandeis University
Josh Falk, University of Chicago
Sarah Fox, Eastern Michigan University
Bethany Greenbaum, Brandeis University
Amy Hemmeter, University of Michigan
Gavriel Hirsch, Northwestern University
Aaron Klein, Harvard University
Andrew Lamont, Eastern Michigan University
Jeffrey Ling, Harvard University
Alexa Little, Yale University
Tom McCoy, Yale University
Yiwei Luo, Princeton University
Jenny Nitishinskaya, Harvard University
Catherine Sheard, University of Oxford
Miriam Wong, Brandeis University
Chelsea Voss, Massachusetts Institute of Technology

Booklet Editors:

Andrew Lamont, Eastern Michigan University
Dragomir Radev, University of Michigan

Sponsorship Chair:

James Pustejovsky, Brandeis University

Corporate, Academic, and Government Sponsors:

The Feldberg Family Foundation
Brandeis University
University of Michigan
Carnegie Mellon University
North American Chapter of the Association for Computational Linguistics
Linguistic Society of America
Mozilla
Association for Computing Machinery, SIGIR
Lockheed Martin
Masco
Yahoo!
Many generous individual donors

Special thanks to:

Tatiana Korelsky, Joan Maling, and D. Terrence Langendoen, US National Science Foundation
And the hosts of the 90+ High School Sites

All material in this booklet © 2014, North American Computational Linguistics Olympiad and the authors of the individual problems. Please do not copy or distribute without permission.



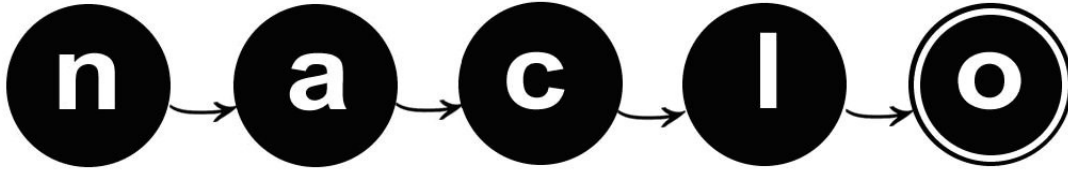
NACLO 2014

Sites



As well as more than 90 high schools throughout the USA and Canada

Please scan the booklet beginning with the next page



The North American Computational Linguistics Olympiad
www.naclo.cs.cmu.edu

Contest Booklet

REGISTRATION NUMBER

--	--	--	--

Name: _____

Contest Site: _____

Site ID: _____

City, State: _____

Grade: _____

Start Time (part I): _____

End Time (part I): _____

Start Time (part II): _____

End Time (part II): _____

Please also make sure to **write your registration number and your name on each page** that you turn in.

SIGN YOUR NAME BELOW TO CONFIRM THAT YOU WILL NOT DISCUSS THESE PROBLEMS WITH ANYONE UNTIL THEY HAVE BEEN OFFICIALLY POSTED ON THE NACLO WEBSITE IN LATE MARCH.

Signature: _____

Part I
Problems I-N
3 Hours

You may only work on this part before the break

YOUR NAME:

REGISTRATION #

(I) To play or not to play (1/2) [10 points]

Kiswahili is a Bantu language with heavy Arabic influence spoken throughout East Africa. While only about 5 million people speak Kiswahili as their first language, over 60 million people use it in their daily life. Kiswahili is an official language of Tanzania, Kenya, Uganda, the Comoros, and the Democratic Republic of the Congo.

11. Match the words in column A with their translations in column B (each translation will be used exactly once):

Column A (Kiswahili)	
1.	Atacheza
2.	Mlifahamu
3.	Mnapika
4.	Nilicheza
5.	Ninapika
6.	Nitapika
7.	Tulifahamu
8.	Unacheza
9.	Utapika
10.	Wanafahamu
11.	Watapika
12.	Walicheza

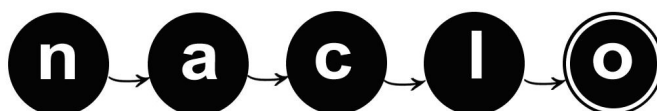
Column B (English)	
A	He/she will play
B	I played
C	I cook
D	I will cook
E	They understand
F	They will cook
G	They played
H	We understood
I	Y'all ¹ understood
J	Y'all cook
K	You play
L	You will cook

12. Match the words in column A with their translations in column B (each translation will be used exactly once):

Column A (Kiswahili)	
1.	Hakucheza
2.	Hamkupika
3.	Hatacheza
4.	Hatapika
5.	Hatukufahamu
6.	Hatupiki
7.	Hawafahamu
8.	Huchezi
9.	Sikucheza

Column B (English)	
A	He/she did not play
B	He/she will not cook
C	He/she will not play
D	I did not play
E	They do not understand
F	We did not understand
G	We do not cook
H	Y'all did not cook
I	You do not play

¹Y'all is the plural form of you



(I) To play or not to play (2/2)

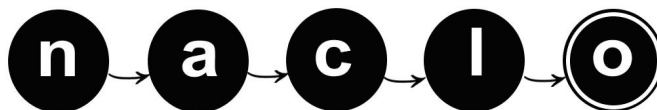
13. Now, here is a class of Kiswahili verbs that work slightly differently. Again, match the words in column A with their translations in column B (each translation will be used exactly once):

Column A (Kiswahili)		Column B (English)	
1.	Hamtakula	A	I do not eat
2.	Hatupi	B	They did not eat
3.	Hawakula	C	They did not give
4.	Hawakupa	D	They will give
5.	Huchi	E	We do not give
6.	Mlikucha	F	Y'all feared
7.	Sili	G	Y'all will not eat
8.	Unakucha	H	You do not fear
9.	Watakupa	I	You fear

14. Given that *ninatembelea* means "I visit" and *ninakufa* means "I die," translate the following into Kiswahili:

1.	You visit.	
2.	I do not visit.	
3.	Y'all visited.	
4.	We did not visit.	
5.	He/she will visit.	
6.	They will not visit.	

7.	You die.	
8.	I do not die.	
9.	Y'all died.	
10.	We did not die.	
11.	He/she will die.	
12.	They will not die.	



YOUR NAME:

REGISTRATION #

(J) Lexicondensed (1/4) [15 points]

Compiling a lexicon (a catalog of words) can be time-consuming and difficult because each individual word has so many potential forms. Suppose that you are dealing with the following words:

view, viewed, viewing, views, review, reviewed, reviewing, reviews, watch, watched, watches, watching, rewatch, rewatches, rewatching, rewatched, wave, waved, waves, waving, rewave, rewaves, rewaved, and rewaving.

Writing all of these forms is tedious; even though you generate a list, you will probably feel listless. Therefore, instead of using this brute force method, you can condense the list with the format shown below:

VERBPREFIX	VERBSTEM	VERBSUFFIX
re	watch	ed
∅	view	s
	wave	ing
		∅

This setup generates a list of all words that consist of one component of VERBPREFIX followed by one component of VERBSTEM followed by one component of VERBSUFFIX (the ∅ stands for an empty spot, so a word could have no letters in the VERBPREFIX or VERBSUFFIX slot). The list generated is identical to the brute force list but is much less tedious to create.

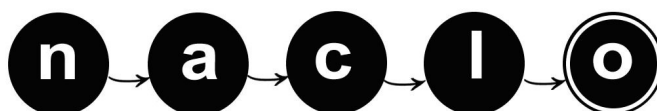
There is one major problem, however. The way that this format strings together word components (called morphemes) does not account for spelling changes that may occur along the way. For example, many legitimate words are generated, such as watch, review, and rewaves, but some misspelled words also result, such as watchs and waving. In order to fix this, you also need to write a set of spelling change rules to describe these changes. The applicable rules in this case are:

```
ch -> che || * s
e -> ∅ || * [ed | ing]
```

These rules mean “ch turns into che if ch is followed by s” and “e turns into nothing if e is followed by ed or ing.”

There are many different ways that this type of rule can be written. Here are a few more examples of such rules and their meanings:

u -> w * Vowel	(u turns into w if u is followed by a vowel)
np -> mp	(np always turns into mp)
t -> c Consonant * kf	(t turns into c if it is between a consonant and kf)
[l f r] -> z w * [c p]	(each letter l, f, or r will turn into z if it falls between w and either c or p)



(J) Lexicondensed (2/4)

J1. Consider the following lexicon and set of rules. (Note that the rules apply in the order given).

PARTONE	PARTTWO	Spelling Change Rules:
cdn	rgt	vsk -> ko
cav	sks	nbj -> jirj
---	---	nsk -> jeej
		gt -> e avr *
		j -> res avb *
		j -> tu b *
		gt -> ar
		vb -> yp
		cdj -> b
		c -> cal q * y
		js -> ch
		os -> o ak *
		ak -> jinkcj c *
		cj -> g
		dnr -> ed
		s -> ry o *
		q -> hi * ck
		q -> eu * ca
		ay -> y * p
		qc -> po
		c -> m * av
		vr -> pl

A. Write the four words generated by the above lexicon and set of rules.

--	--	--	--

B. If you add two more three-letter entries to the lexicon (one entry in PARTONE and one entry in PARTTWO), the system will generate an additional five words that go together with the four words from Task I. What are the new entries for PARTONE and PARTTWO?

--	--

What are the five newly generated words? (Hint: Every rule is used at least once.)

--	--	--	--	--



YOUR NAME:

REGISTRATION #

(J) Lexicondensed (4/4)

Country	Desired Adjective	Country	Desired Adjective
andorra	andorran	japan	japanese
australia	australian	kenya	kenyan
bhutan	bhutanese	mexico	mexican
bolivia	bolivian	morocco	moroccan
cambodia	cambodian	nauru	nauruan
chad	chadian	netherlands	dutch
chile	chilean	poland	polish
china	chinese	portugal	portuguese
congo	congolese	rwanda	rwandan
cuba	cuban	singapore	singaporean
cyprus	cyriot	sudan	sudanese
england	english	togo	togolese
fiji	fijian	uganda	ugandan
guyana	guyanese	vietnam	vietnamese
indonesia	indonesian	yemen	yemeni
israel	israeli		



YOUR NAME:

REGISTRATION #

(K) Don't be Ukhrl to a Liver that's True (1/2) [10 points]

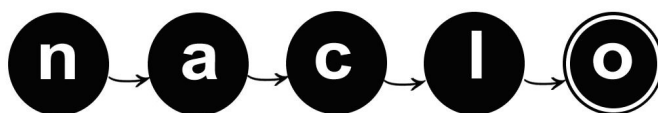
As you may know, languages form “families” in which languages descended from a common ancestor (ancient language) show systematic similarities and differences. For example English, Dutch and Danish are all from the same language family, and the systematic difference can be seen in the words for *brother*, *mother*, *father* in Dutch (*broeder*, *moeder*, *vader*) and Danish (*bror*, *mor*, *far*). The French words *frère*, *mère* and *père* are also (more distantly) related, and show slightly more complex differences.

Kachai, Tusom, and Ukhrl are three languages from the Tangkhulic subfamily of the Tibeto-Burman family of languages. They are spoken in Manipur state, India. The words from these languages that are given here form sets of three that are descended from the same word in the shared ancestor of the three languages. The Ukhrl words are given in the table on the next page, with their English translations. Kachai and Tusom words are given in no particular order. Write the letters corresponding to the Kachai and Tusom words in proper order in the table on the following page.

Pronunciation notes:

- The small raised *h* symbol indicates that the preceding consonant is aspirated, i.e. pronounced with an exaggerated puff of air.
- ə represents a vowel like the first sound of the word *approach*.
- ŋ represents a velar nasal, the ‘ng’ sound in a word like *sing*.
- ʔ is a glottal stop, the sound between the two syllables of the expression *uh-oh*.
- ð is the ‘th’ sound at the beginning of *this*.
- ɐ represents a vowel somewhere between the ‘a’ in *cat* and the ‘o’ in *cot*.
- x is pronounced like ‘ch’ in *Bach*.
- ʊ is a vowel pronounced like ‘oo’ in *book*, but with spread lips, a bit like when you show distaste *ugh*
- ɔ̃ is nasal vowel, pronounced like the ‘on’ in *bon vivant*.
- ʃ is the ‘sh’ sound in *ship*
- c is pronounced like ‘ch’ in *church*

Kachai				Tusom			
(A)	k ^h əŋət ^h i	(K)	kəkwe	(a)	kət ^h ue	(k)	ʃi
(B)	kəp ^h u	(L)	ʔami	(b)	kətχa	(l)	ma
(C)	mək ^h u	(M)	ʔamɐ	(c)	mu	(m)	mokʃi
(D)	ʔamət ^h ɐn	(N)	kəce	(d)	kəkie	(n)	luə
(E)	ʔale	(O)	ʔacu	(e)	k ^h əŋie	(o)	ʔətχa
(F)	k ^h əmwe	(P)	kət ^h e	(f)	ʔəntsɔ̃	(p)	za
(G)	ʔat ^h i	(Q)	k ^h əmɐn	(g)	k ^h anny	(q)	ci
(H)	kək ^h u	(R)	kət ^h i	(h)	k ^h antsy	(r)	k ^h əmɔ̃
(I)	kəði	(S)	ʔak ^h we	(i)	kʃie	(s)	makəcuə
(J)	ʔasu	(T)	k ^h əməni	(j)	kəpʃi	(t)	kəkʃi



YOUR NAME:

REGISTRATION #

(K) Don't be Ukhru! to a Liver that's True (2/2)

Kachai	Tusom	Ukhru!	English
		kət ^h uj	awaken
		kək ^h a	bitter
		kəkaj	break
		kəcu!	burn
		k ^h əŋaj	desire
		k ^h əŋət ^h u	exchange
		lu!	field
		me!	fire
		sa	flesh/animal
		ʔat ^h ej	fruit
		mi	human
		mək ^h a	jaw
		k ^h aj	knife
		k ^h əmənu	laugh
		ʔamət ^h in	liver
		ca	necklace
		k ^h əmin	ripe
		kət ^h ej	see
		kəp ^h a	seek
		tse!	spear



YOUR NAME:

REGISTRATION #

(L) Transducing Runes (1/5) [10 points]

Before the Roman alphabet was introduced to Northern Europe, much of Scandinavia and what is now Great Britain used a writing system called Runic. These symbols have recently gained increasing popularity because the fantasy author J.R.R. Tolkien adapted an Anglo-Saxon Runic writing called Futhorc in his series *Lord of the Rings* (and *The Hobbit*).

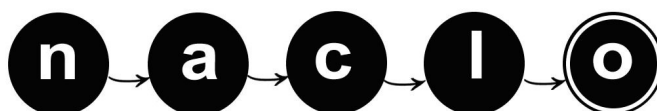
This problem is about mathematical constructs that we can use to turn Roman text (i.e., what English is written in) into runes. This is not a simple substitution, however, because there is not a one-to-one connection between Roman letters and runes. For example, these words become the following runes. To make things cleaner, we're assuming that every word written in Roman characters is followed by a # to mark the end of the word. You can assume that every input Latin word will be terminated by a #, and that this becomes ■ in runes.

Roman	Runic
sat#	↵ ↶ ↑ ■
eat#	↵ ↑ ■
heat#	↵ ↶ ↵ ↑ ■
east#	↵ □ ■

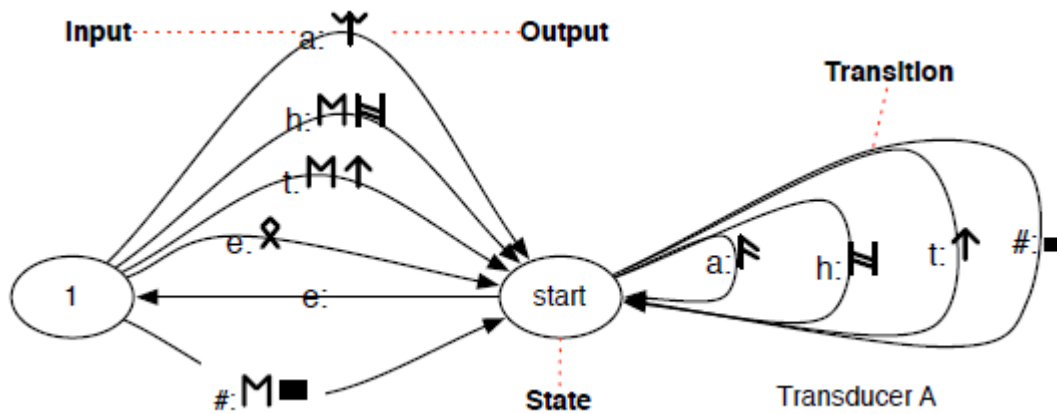
Specifically, there are a number of runes that are equivalent to two Roman characters. To keep things simple, we'll start with a very limited alphabet.

a	↶	ea	↵↶
e	↶↶	ee	↶↶
h	↶↶↶	th	↶↶
s	↶↶↶↶	st	↶↶↶
t	↑	#	■

The tool that we're going to use is called a **transducer**, a logical tool that is used in morphological processing (e.g., to remove suffixes and prefixes from words) in natural language processing technology.



(L) Transducing Runes (2/5)



The key components of a transducer are states, transitions, inputs, and outputs. We always start in the “start” state. In the example transducer below, this is the right circle with the label “start” inside it. We transition to different states based on the input that we get.

In this problem, our input is Roman characters. For example, if we're in the “start” state and see either *h*, *a*, or *t*, then we transition from the “start” state to the “start” state (simple!). If, however, we were in the “start” state and saw the character *e*, we would transition to state “1”.

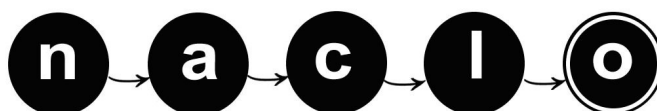
Which transition we use is based on the input we receive. When we transition, we also can output. In the start state,

- if we see *h* we output 𐌆;
- if we see *a* we output ǁ;
- if we see *t* we output 𐌇;
- if we see *#* we output 𐌹;
- but if we see *e* we output nothing.

Transitions are depicted with an arrow. Each arrow has a label that shows the input and output. To the left of the colon (:) is the input, and the output is to the right (possibly empty, as in the case of *e* in the start state).

Different states can have different transitions; we output different runes based on input. In state “1”; for example, if we then see *a*, we output ǀ, which allows us to turn the input of *e* followed by *a* into the correct rune. Thus, if we're in state “1” it means that we might need to turn a sequence of characters into a single rune, but we won't know for sure until we see the next character.

If you're unclear on the concept, trace *eat#* and *heat#* through this simple transducer and make sure you get outputs that match the example runes.

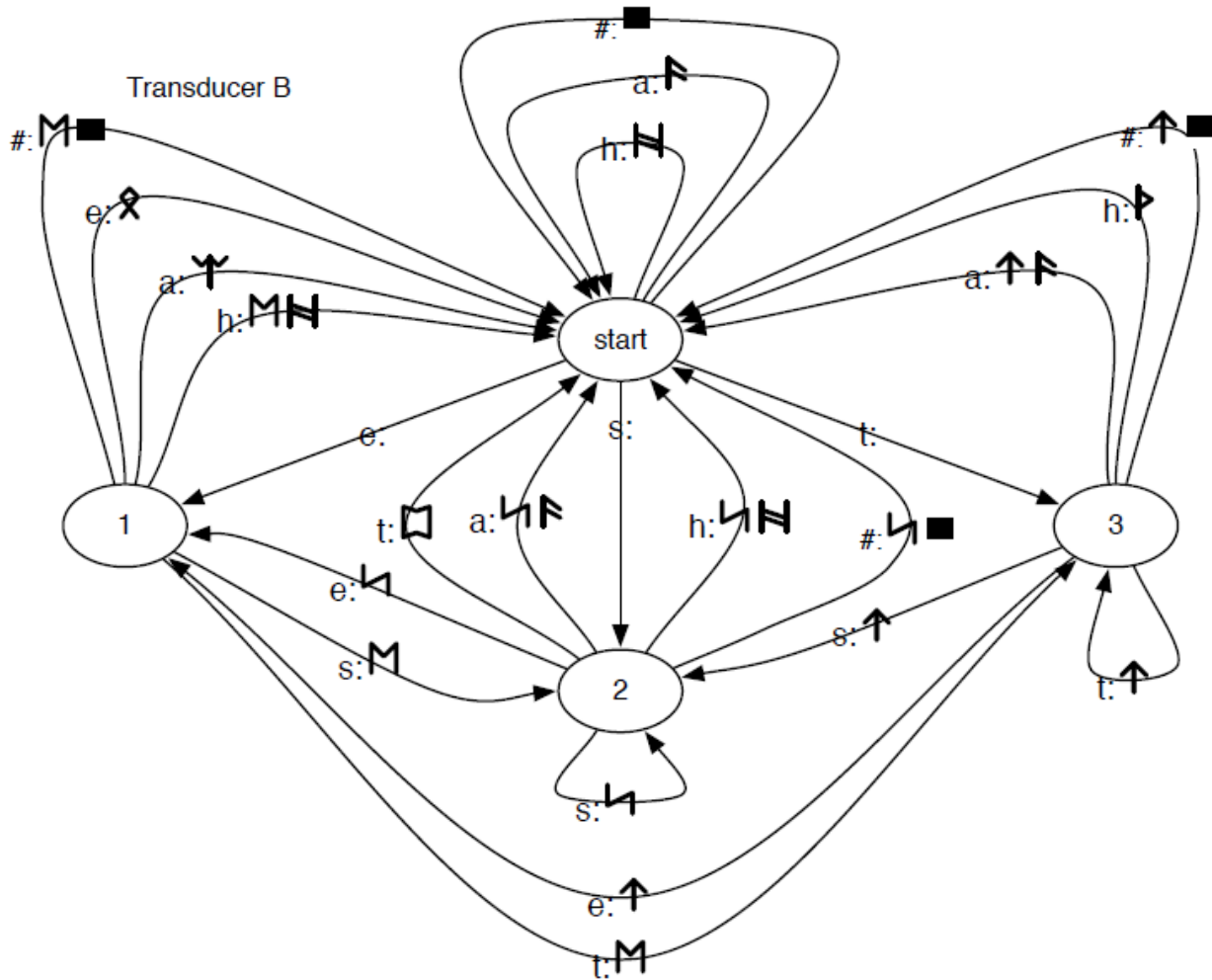


YOUR NAME:

REGISTRATION #

(L) Transducing Runes (3/5)

L1. Below is a transducer for the letters *a*, *e*, *h*, *s*, *t*, and *#*. Given a sequence of Roman characters, give the states that you would visit while transducing those characters. The first is done as an example.



A)	he#	start	start	l	start			
B)	stash#	start						
C)	heath#	start						
D)	thee#	start						

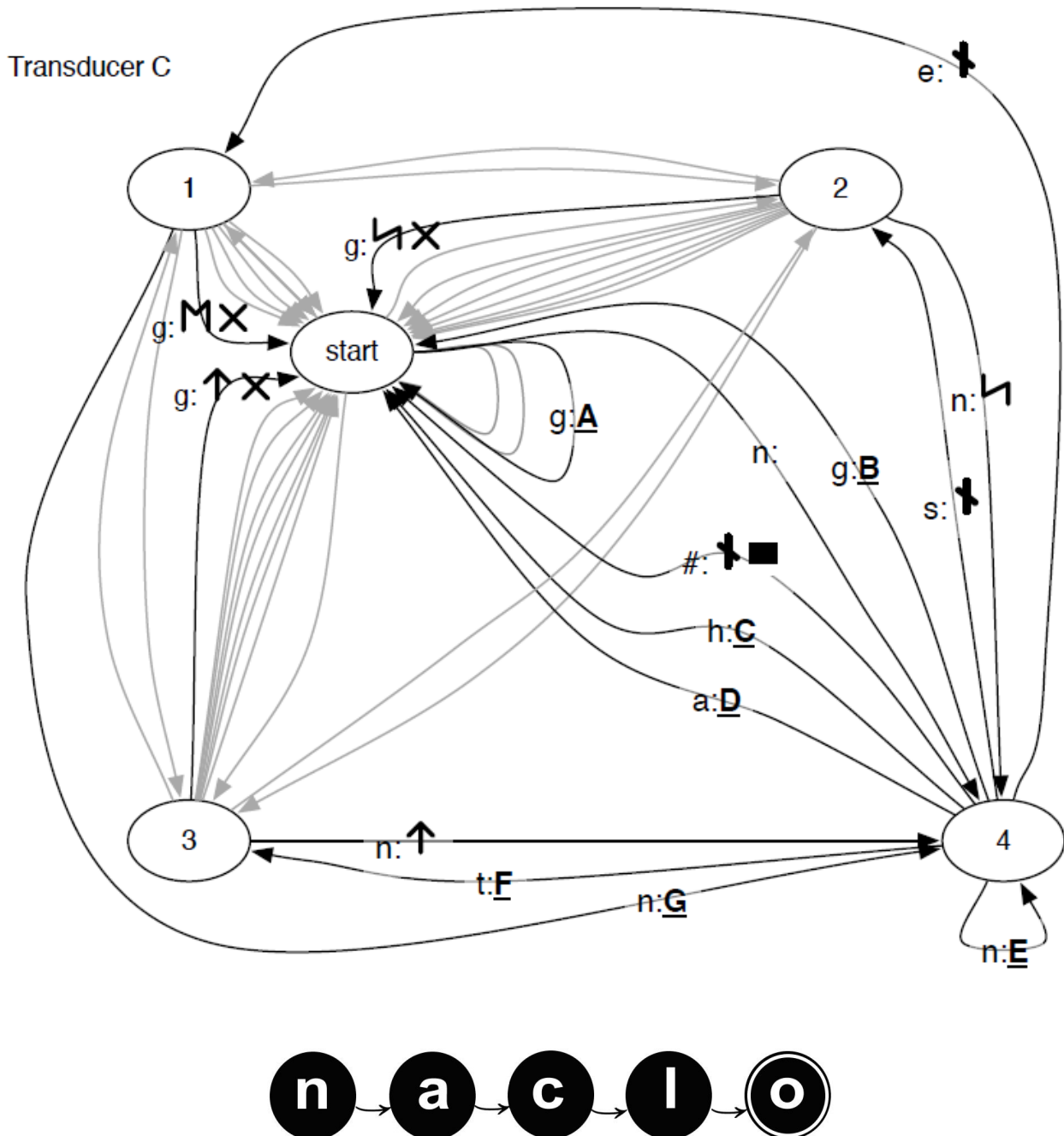


(L) Transducing Runes (4/5)

L2. We're going to make our transducer a little more complicated, by adding additional runes. The additional runes we'll add correspond to the letters *n*, *g*, and *ng*.

n † *g* ✕ *ng* ✕

Below is what this transducer looks like. It's getting more complex, so we're not going to show all of it. Instead, we'll show transitions that were in the previous transducer in gray without the inputs and outputs. We also won't give the outputs for some of the transitions; some of the outputs have been replaced by bold, upper-case, underlined Roman letters; you'll fill in those missing runes on the next page.



(L) Transducing Runes (5/5)

What is the correct output for the transitions in the above transducer? Use the numbered runes below. CAUTION: Answers can be repeated, outputs may require more than one rune, and *order matters*.

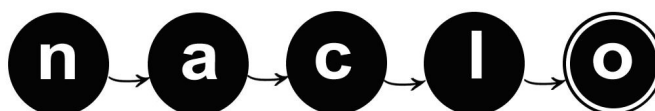
- | | | | |
|------------------|-------|------------------|-------|
| (a) g : <u>A</u> | _____ | (e) n : <u>E</u> | _____ |
| (b) g : <u>B</u> | _____ | (f) t : <u>G</u> | _____ |
| (c) h : <u>C</u> | _____ | (g) n : <u>H</u> | _____ |
| (d) a : <u>D</u> | _____ | | |

1. **X** 2. **†** 3. **M** 4. **N** 5. **F** 6. **L** 7. **X** 8. **↑**

L3. Consider the number of states and transitions in a transducer needed to represent different alphabets. The table has the number of states and transitions for the transducers previously shown (don't forget about the end of the word marked #).

Transducer	Single Runes	Double Runes	States	Transitions
A	F (a), N (h), ↑ (t)	M (e), ↑ (t)	2	10
B	F (a), N (h), ↑ (t)	M (e), L (s), ↑ (t)	4	24
C	F (a), N (h), † (n), ↑ (t)	M (e), X (g), L (s), ↑ (t)	5	?
D	F (a), M (e), X (g), L (s),	X (d), N (h), † (n), ↑ (t)	?	?

A)	How many transitions does transducer C have?	
B)	How many states does transducer D have?	
C)	How many transitions does transducer D have?	



YOUR NAME:

REGISTRATION #

(M) Come to Istanbul (I/I) [10 points]

Turkish is spoken by about 63 million people, of whom most live in Turkey but about 100,000 live in the UK. It is a non-Indo-European language, so it is unrelated to English but related to languages of Central Asia such as Azeri and Uzbek.

Turkish words are built up by adding one or more endings to a root word; the vowels in most word endings vary depending on the vowels in the root word ("vowel harmony"), as you will see in the following examples. Here are some sentences in Turkish, with their English translations. Note:

- The Turkish letters "ş", "ç" and "ı" are pronounced like English "sh", "ch" and the "a" in "above".
- The letters i and ı represent different vowels.
- The letter "ğ" is usually silent (like the "gh" in "although").
- Square brackets [] enclose English words that are not directly translated.

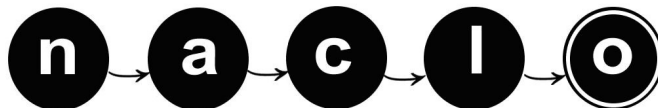
Arkadaşlarım şehirde mutlu	My friends [are] happy in [the] city.
Baban İstanbul'u seviyor mu?	Does your father like Istanbul?
Fakirler Van'dan İstanbul'a gelmek istiyor	Poor [people] want to come from [the city of] Van to Istanbul.
İstanbul en büyük şehir	Istanbul [is the] biggest city.
Eve geliyorlar	They come home.
Babam "Merhaba! Gel, arkadaşımız ol", diyor	My father says "Hello! Come [and] be our friend".
Evimizde büyük pencereler var	There are big windows in our house.
Pencereden atlıyoruz	We jump from [the] window.
Ev almak mı istiyorsun?	Do you want to buy [a] house?

M1. How would you translate the following into English?

A.	Baban mutlu mu?	
B.	"Şehrimize gel" diyoruz.	
C.	Arkadaşım doktor olmak istiyor.	
D.	Fakir evimi seviyorlar mı?	
E.	İstanbul'dan mı geliyorsun?	

M2. The following examples introduce a new pattern. What do you think these examples mean?

A.	Geldiğimde "merhaba" diyorlar.	
B.	Baban geldiğimizden mutlu mu?	
C.	Fakir olduğunu diyorlar.	
D.	Aldığın ev büyük mü?	
E.	En mutlu olduğum şehir, Van.	
F.	Fakir olduğumuz halde mutluyuz.	



YOUR NAME:

REGISTRATION #

(N) Hungarian Rocks (I/I) [5 points]

The grid below represents a field divided into a 7 x 7 grid, aligned north-south and east-west. In some of the squares of the grid are rocks represented by X.

There are four Hungarians – Dorottya, László, Erika, and Balázs – standing in the field, each in a different square not containing a rock, and each facing in one of the four cardinal directions (north, south, east west) - not necessarily different from each other. Each person makes some statements describing the positions of the rocks. For instance, Dorottya’s first statement means “(Due) east (behind me) there is one stone.”

Find each person’s place in the field and the direction they are facing. References to directions are to be understood as describing a single line in the field: “due east”, “directly behind me”, and so on.

Dorottya says: Keletre (mögöttem) egy kő van.
 Délre két kő van.
 Jobbra nincs kő.

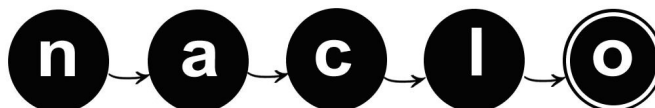
László says: Délre (balra) nincs kő.
 Északra egy kő van.
 Mögöttem két kő van.

Erika says: Északra (előttem) nincs kő.
 Nyugatra egy kő van.
 Jobbra két kő van.

Balázs says: Nyugatra (jobbra) két kő van.
 Északra egy kő van.
 Balra nincs kő.

Position	Direction

	A	B	C	D	E	F	G
1				X			
2			X	X	X		
3							
4	X			X			X
5			X		X		
6							
7				X			



Part 2
Problems O-Q
2 Hours

You may only work on this part after the break

YOUR NAME:

REGISTRATION #

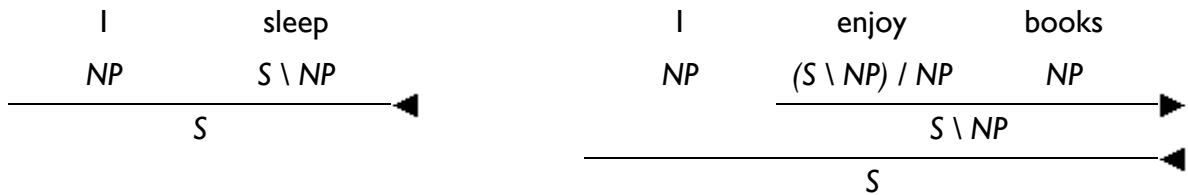
(O) CCG (1/2) [5 points]

One way for computers to understand language is by forming a structure that represents the relationships between words using a technique called Combinatorial Categorical Grammar (CCG). Computer scientists and linguists can use CCG to parse sentences (that is, try to figure out their structure) and then extract meaning from the structure.

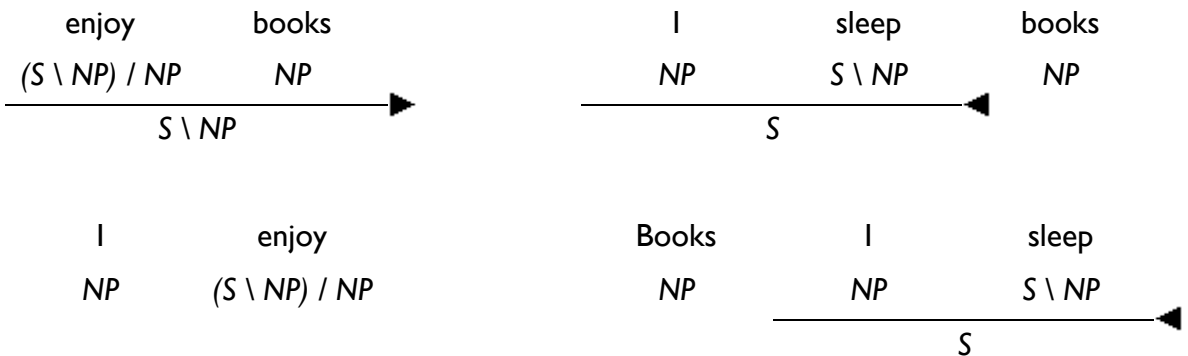
As the name suggests, Combinatorial Categorical Grammar parses sentences by combining categories. Each word in a sentence is assigned a particular category; note that / and \ are two different symbols:

I	NP
books	NP
sleep	$S \setminus NP$
enjoy	$(S \setminus NP) / NP$

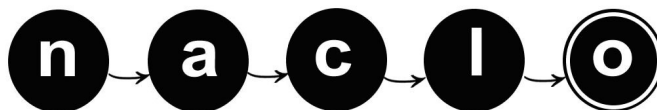
These categories are then combined in systematic ways. We will not explain how, but we will give you two successful parses...



...and four unsuccessful parses...



If a parse is successful, the sentence is declared “grammatical”; if not, the sentence is declared “ungrammatical”.



YOUR NAME:

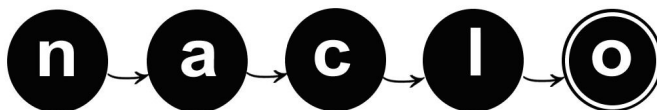
REGISTRATION #

(O) CCG (2/2)

O1. Using the above examples as evidence, figure out how CCG parses sentences, and describe it briefly here:

O2. In the sentence “I enjoy long books”, list all of the categories that, if assigned to “long”, make the sentence have a successful parse.

O3. Not every grammatical sentence of English will be declared “grammatical” by the process above. Using only the words “I”, “books”, “sleep”, and “enjoy”, form a grammatically correct English sentence that will fail to parse given the categories above. You don’t have to use all four of the words.



YOUR NAME:

REGISTRATION #

(P) Combining Categories in Tok Pisin (1/2) [15 points]

This problem is a follow-up to problem O and has to be solved after that problem. Tok Pisin (also referred to as New Guinea Pidgin or Melanesian Pidgin) is a creole language spoken in the northern mainland of Papua New Guinea and surrounding islands. It is an official language and the mostly widely used language in the country, spoken by over 5 million people.

Many Tok Pisin words come originally from English – its name comes from “talk” and “pidgin”¹ -- but Tok Pisin isn't just English. It has a distinct grammar and uses these words in different (but systematic!) ways.

P1. Below are sentences in Tok Pisin with a scrambled list of English translations. Match each sentence to its English equivalent.

1.	Brata bilong em i stap rit.	
2.	Ol i stap dringim wara.	
3.	Ol i ken ritim buk bilong mi.	
4.	Em i ritim buk pinis.	
5.	Em i laik rit.	
6.	Susa bilong em i ken rait.	
7.	Susa bilong mi i boylim wara.	
8.	Wara i boil pinis.	

A.	He has read the book.
B.	My sister boils the water.
C.	They can read my book.
D.	His sister can write.
E.	His brother is reading.
F.	The water has boiled.
G.	He wants to read.
H.	They are drinking water.

P2. Translate the following Tok Pisin sentence into English:

Brata bilong mi i stap ritim buk bilong susa bilong mi.

P3. Translate the following English sentence into Tok Pisin:

Their sister wants to write a book.

¹A pidgin language is a communicative system developed by two or more groups of people who do not share a common language. Tok Pisin started out as a pidgin but has since developed into a creole, a complex language in its own right.



YOUR NAME:

REGISTRATION #

(P) Combining Categories in Tok Pisin (2/2)

P4. Describing these words in terms of their CCG categories (introduced in Problem O) highlights that these aren't English words combined according to English rules, but are Tok Pisin words combined according to Tok Pisin rules.

Match each Tok Pisin word to its CCG category. Some categories will be used more than once. The symbol S_b is short for 'Bare Clause'.

1.	bilong	
2.	brata	
3.	boil	
4.	boilim	
5.	buk	
6.	dringim	
7.	em	
8.	i	
9.	ken	
10.	laik	

11.	mi	
12.	ol	
13.	pinis	
14.	stap	
15.	raitim	
16.	rit	
17.	ritim	
18.	susa	
19.	wara	

A.	NP
B.	$(NP \setminus NP) / NP$
C.	$(S \setminus NP) / (S_b \setminus NP)$
D.	$(S_b \setminus NP)$
E.	$(S_b \setminus NP) / NP$
F.	$(S_b \setminus NP) \setminus (S_b \setminus NP)$
G.	$(S_b \setminus NP) / (S_b \setminus NP)$

P5. Explain your answer.



YOUR NAME:

REGISTRATION #

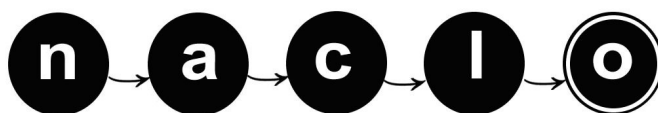
(Q) Learning Yidiny (1/2) [20 points]

Yidiny is the language of people whose ancestral lands are in the rain forest country of northeastern Queensland, Australia, south of Cairns. Here are some Yidiny sentences recorded from mother tongue (or first language) speakers of this language.¹

Examine sentences (1) to (21) and try to work out the meaning of each word and why words with the same meaning may have different forms. Sometimes a single word of Yidiny may need to be translated by two – or even several – English words; the converse may also be true. The given translations are in order.

- | | |
|--|--|
| 1. Nganji jarral dunggul guluguluugu. | <i>We set up a fish-trap for black bream.</i> |
| 2. Nganjiiny bamaal gugaal mayiigu | <i>The people called us for food.</i> |
| 3. Wanjiirr nyuniinda mayi? | <i>How much food have you got?</i> |
| 4. Ngayu banjaar gabay. | <i>I followed the road.</i> |
| 5. Ngayu biwuuda minya jaban bagaal. | <i>I speared an eel with a fish-spear.</i> |
| 6. Nganji dugur balgaal jirrgaada. | <i>We made a hut with grass.</i> |
| 7. Nganyany jina banggaaldu gundaajinyu. | <i>The axe happened to cut my foot.</i> |
| 8. Ngayu waguuja banggaalda gundaal. | <i>I cut the man with an axe.</i> |
| 9. Nganyany wagujanggu banggaalda gundaal. | <i>The man cut me with an axe.</i> |
| 10. Nyundu gana nganda guman wiwin. | <i>You just give me one.</i> |
| 11. Ngayu nyuniny wawaal. | <i>I saw you.</i> |
| 12. Nganyany bamaal wawaal. | <i>A person saw me.</i> |
| 13. Ngayu bama wawaajinyu jambuul. | <i>I happened to see two people.</i> |
| 14. Minyaagu yingu gadang jabaangu. | <i>This (one) is coming for eels.</i> |
| 15. Ngayu bama bunya barrgandanyu. | <i>I passed the woman by.</i> |
| 16. Nganyany bamaal bunyaang barrgandanyu. | <i>The woman passed me by.</i> |
| 17. Ngungu bunya gabaanja janaany. | <i>That woman was standing on the road.</i> |
| 18. Nganjiinda jaja ngunjuung ngurrangurraal bunyaang. | <i>That woman showed us the baby.</i> |
| 19. Waguuja dungu bunyaang jinaa baraal. | <i>The woman kicked the man in the head.</i> |
| 20. Bunya wagujanda dunguu jinaa baraajinyu. | <i>The woman happened to kick the man in the head.</i> |
| 21. Ngayu bama mandii baraal. | <i>I punched the person.</i> |

¹Yidiny was described by linguist RMW Dixon in his 1977 book entitled *A grammar of Yidiny* published by Cambridge University Press. Sentences (1-21) are from this publication, with the original IPA (phonetic) symbols transliterated.



(Q) Learning Yidiny (2/2)

The sentences (A) to (L) below were spoken by a person who is not a native speaker of Yidiny, who was trying to learn Yidiny as a second language. This speaker makes grammatical mistakes. The English sentences indicate what the speaker was trying to say in Yidiny. In each of these ungrammatical sentences (indicated by the asterisk (*)) an incorrect form of *one* word is used. Your task is to locate the ungrammatical word in each sentence. Copy it into the appropriate column of the table below, and then write the correct form of the word in the column to the right of the incorrect word form. (Don't worry about the Yidiny word order.)

- | | |
|---|--|
| A. *Nyuniny gabay mijil. | <i>You are blocking the road.</i> |
| B. *Ngayu nyuniny jina banggaaldu gundaal. | <i>I cut your foot with an axe.</i> |
| C. *Nganjiiny bama bunyaang wawaal. | <i>The woman saw us.</i> |
| D. *Wanjiirr ngayu minya? | <i>How much meat do I have?</i> |
| E. *Bamaal waguuja gabaanja janaany. | <i>The man was standing on the road.</i> |
| F. *Nganji ngungu guluguluugu bagaal. | <i>We speared that black bream.</i> |
| G. *Bama ngungu dugur balgaal gabaanja. | <i>A person made that hut near the road.</i> |
| H. *Nganjiiny ngungu mayi wiwin. | <i>Give us that food.</i> |
| I. *Nyundu bama bunya mandi bagaal biwuudu. | <i>You stabbed the woman's hand with a fishing-spear.</i> |
| J. *Nyundu jina bagaajinyu biwuudu. | <i>You happened to get stabbed in the foot by a fishing-spear.</i> |
| K. *Nganji jaja dunguu wawaal. | <i>We saw the child's head.</i> |
| L. *Ngayu ngungu bunyaang mandii baraal. | <i>I punched that woman.</i> |

Sentence	Incorrect Word	Corrected Word
A.		
B.		
C.		
D.		
E.		
F.		
G.		
H.		
I.		
J.		
K.		
L.		



YOUR NAME:

REGISTRATION #

Extra Page - Enter the Problem Name Here: _____

