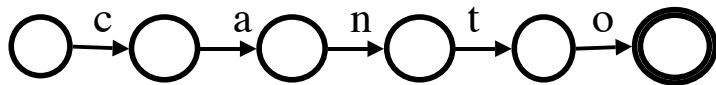


# FST Morphology

Based on Beesley and Karttunen 2003

Miriam Butt  
October 2003

## A One-Word Language



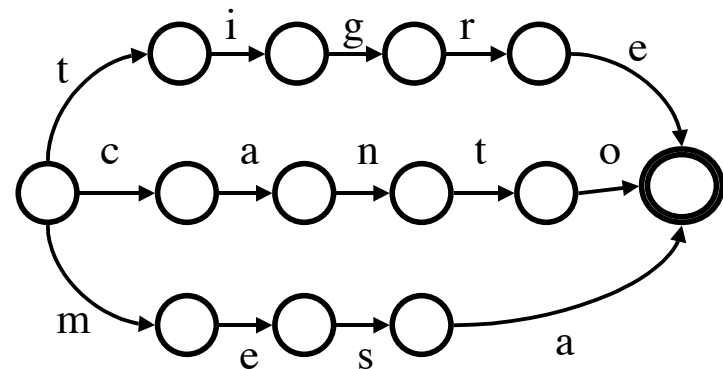
## Recap

**Last Time:** Finite State Automata can model most anything that involves a finite amount of states.

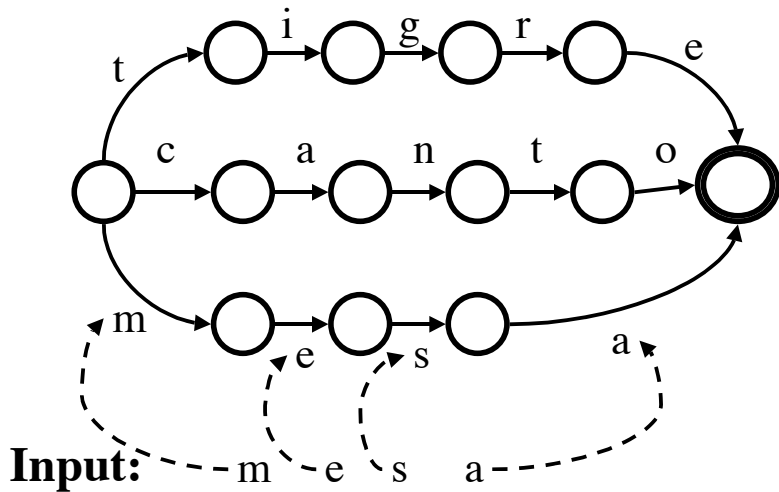
We modeled a Coke Machine and saw that it could also be thought of as defining a *language*.

We will now look at the extension to natural language more closely.

## A Three-Word Language



## Analysis: A Successful Match



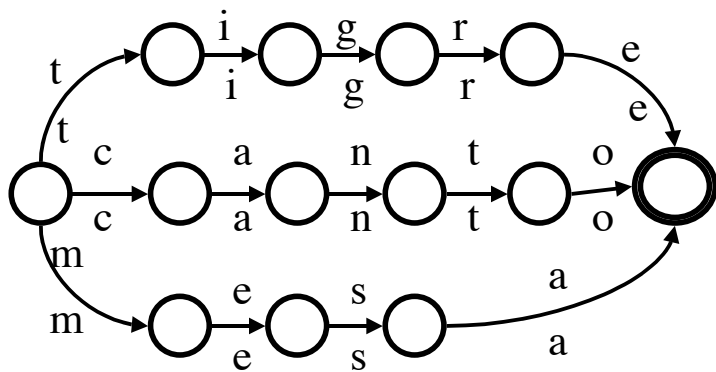
## Rejects

The analysis of *libro*, *tigra*, *cant*, *mesas* will fail.

Why?

**Use:** for example for a spell checker

## Transducers: Beyond Accept and Reject

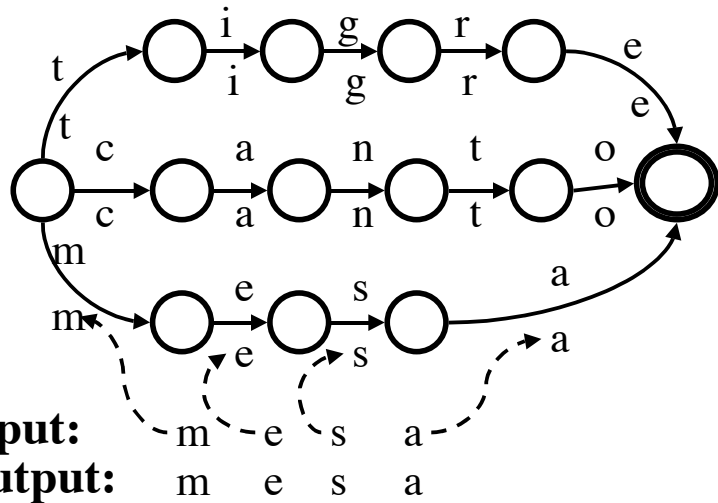


## Transducers: Beyond Accept and Reject

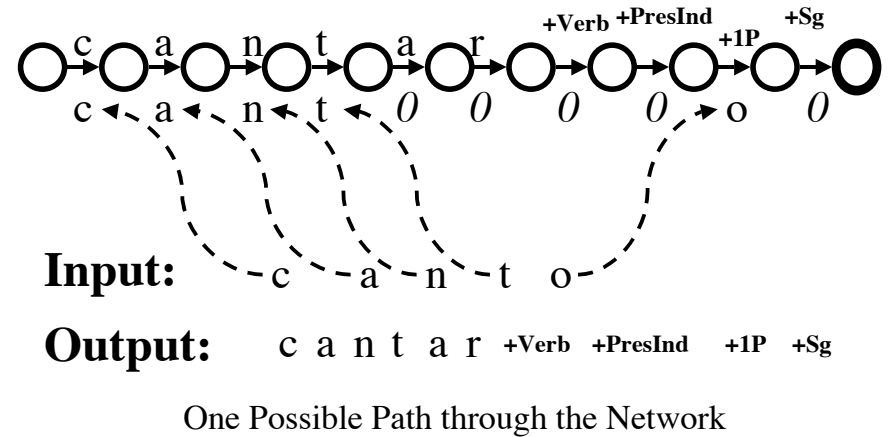
### Analysis Process:

- Start at the Start State
- Match the input symbols of string against the *lower-side* symbol on the arcs, consuming the input symbols and finding a path to a final state.
- If successful, return the string of *upper-side* symbols on the path as the result.
- If unsuccessful, return nothing.

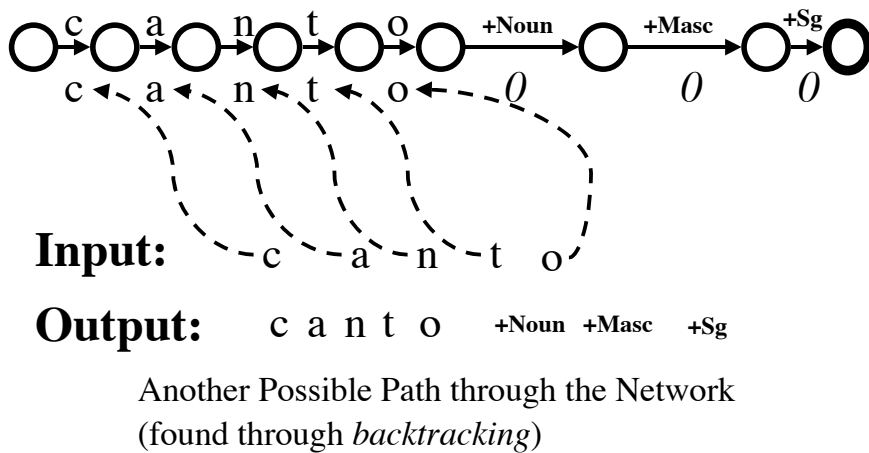
### A Two-Level Transducer



### A Lexical Transducer



### A Lexical Transducer



### Why Transducer?

**General Definition:** Device that converts energy from one form to another.

**In this Context:** Device that converts one string of symbols into another another string of symbols.

## The Tags

Tags or Symbols like +Noun or +Verb are *arbitrary*: the naming convention is determined by the (computational) linguist and depends on the larger picture (type of theory/type of application).

One very successful tagging/naming convention is the Penn Treebank Tag Set

## The Tags

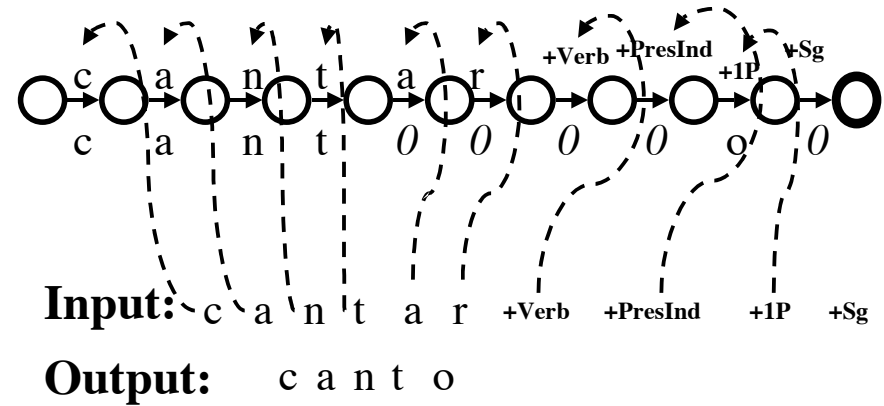
What kind of Tags might be useful?

## Generation vs. Analysis

The same finite state transducers we have been using for the *analysis* of a given surface string can also be used in reverse: for *generation*.

The XRCE people think of analysis as *lookup*, of generation of *lookdown*.

## Generation --- Lookdown



## Generation --- Lookdown

### Analysis Process:

- Start at the Start State and the beginning of the input string
  - Match the input symbols of string against the *upper-side* symbols on the arcs, consuming the input symbols and finding a path to a final state.
  - If successful, return the string of *lower-side* symbols on the path as the result.
  - If generation is unsuccessful, return nothing.

## Sharing Structure and Sets

Networks can be compressed quite cleverly (p. 16-17).

### Sets:

FST Networks are based on formal language theory. This includes basics of set theory:

membership, union, intersection,  
subtraction, complementation

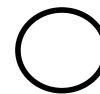
## Tokenization

### General String Conversion: Tokenizer

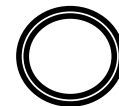
**Task:** Divide up running text into individual tokens

- couldn't -> could not
- to and from -> to-and-fro
- ,, -> ,
- The farmer walked -> The^TBfarmer^TBwalked

## Some Basic Concepts/Representations

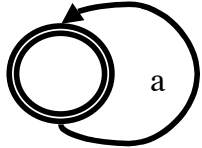


Empty Language

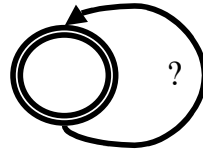


Empty-String Language

## Some Basic Concepts/Representations



A Network for an Infinite Language



A Network for the Universal Language

## Relations

**Ordered Set:** members are ordered.

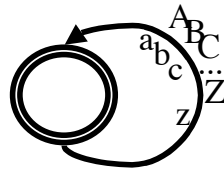
**Ordered Pair:**  $\langle A, B \rangle$  vs.  $\langle B, A \rangle$

**Relation:** set whose members are ordered pairs

- Family Trees (p. 21)
- {  $\langle \text{“cantar+Verb+PresInd+1P+Sg”}, \text{“canto”} \rangle$ ,  $\langle \text{“canto+Noun+Masc+Sg”}, \text{“canto”} \rangle$ , ... }

## Relations

**An Infinite Relation:**



**Identity Relation:**  $\langle \text{“canto”}, \text{“canto”} \rangle$

## Basic Set Operations

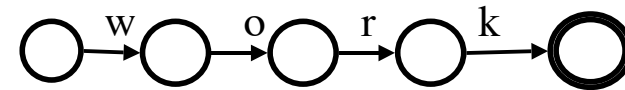
- union (p. 24)
- intersection (p. 25)
- subtraction (p. 25)

## Concatenation

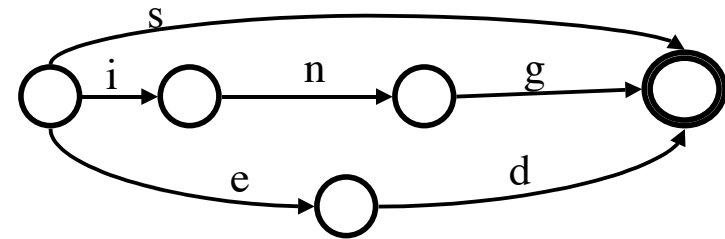
One can also concatenate two existing languages (finite state networks with one another to build up new words productively/dynamically).

This works nicely, but one has to write extra rules to avoid things like: *\*trys*, *\*tryed*, though *trying* is okay.

## Concatenation

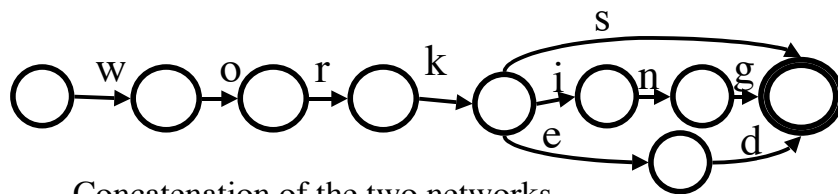


Network for the Language {"work"}



Network for the Language {"s", "ed", "ing"}

## Concatenation



Concatenation of the two networks

What strings/language does this result in?

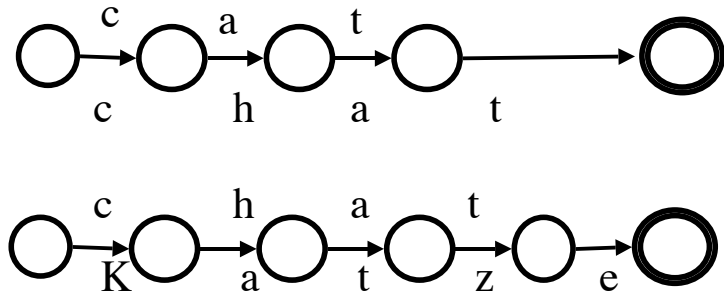
## Composition

Composition is an operation on two relations.

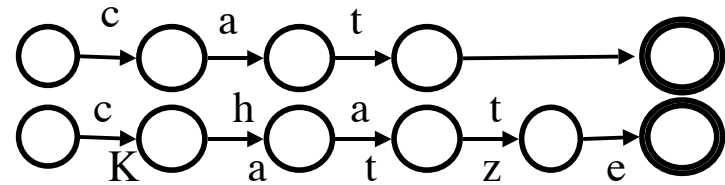
Composition of the two relations  $\langle x, y \rangle$  and  $\langle y, z \rangle$  yields  $\langle x, z \rangle$

Example:  $\langle \text{"cat"}, \text{"chat"} \rangle$  with  $\langle \text{"chat"}, \text{"Katze"} \rangle$  gives  $\langle \text{"cat"}, \text{"Katze"} \rangle$

## Composition

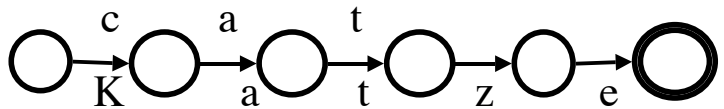


## Composition



Merging the two networks

## Composition



The Composition of the Networks

What is this reminiscent of?

## Composition + Rule Application

Composition can be used to encode phonology-style rules.

E.g., Lexical Phonology assumes several iterations/levels at which different kinds of rules apply.

This can be modeled in terms of cascades of FST transducers (p. 35).

These can then be composed together into one single transducer.



## Next Time

- Begin working with xfst
- **Now:** Exercises 1.10.1, 1.10.2, 1.10.4  
Optional: 1.10.3