# Natural Language Processing

## Timo Baumann

# This week

- Your questions after the first week
- recap: Chomsky hierarchy of formal languages
- our first real topic:
  finite state methods applied to morphology
- exercise: work with finite state transducers (Friday)
  - by hand and using pre-existing software

# Your questions

# Formal Languages and Automata

- Automata theory and theory of formal languages are a part of theoretical computer science

- Their concepts originate in theoretical linguistics: Noam Chomsky is the originator of the Chomsky hierarchy of formal languages

Why talk about it?

- complexity of sub-systems of natural language informs complexity of automatic processing machinery

- Fundamental results from theoretical computer science have direct implications on implementations for language technology applications

# Definitions

- a **LANGUAGE** is a collection of sentences of finite length all constructed from a finite alphabet of symbols

- a **GRAMMAR** can be regarded as a device that enumerates the sentences of a language

- a grammar of language L can be regarded as a function whose range is exactly L

# Formal Grammar

A **formal grammar** is a quad-tuple *G = (Φ,Σ,R,S)* where

- Φ is a finite set of non-terminals

- Σ a finite set of terminals, disjoint from Φ

- R a finite set of production rules of the form
  $$\alpha \in (\Phi \cup \Sigma)^* \rightarrow \beta \in (\Phi \cup \Sigma)^* \text{ with } \alpha \neq \varepsilon \text{ and } \alpha \notin \Sigma^*$$

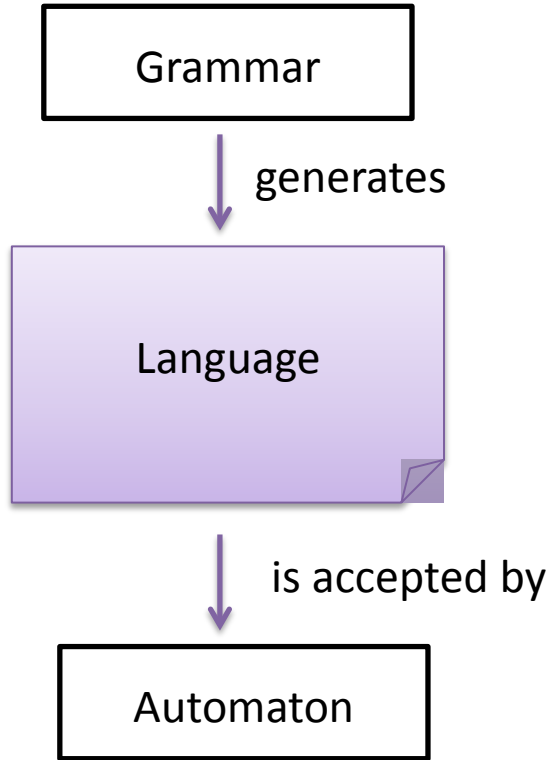- S, Element of Φ : start symbol

# Derivation, Formal Language, Automaton

Let $G = (\Phi, \Sigma, R, S)$ be a formal grammar and let $u, v \in (\Phi \cup \Sigma)^*$.

1. $v$ is **directly derivable** from $u$, noted $u \Rightarrow v$, if
   $u = awb$, $v = azb$ and $w \rightarrow z$ is a production rule in $R$.

2. $v$ is **derivable** from $u$, noted $u \overset{*}{\Rightarrow} v$, if there are words $w_0 .. w_k$, such that $u \Rightarrow w_0$, $w_{n-1} \Rightarrow w_n$ for all $0 < n \leq k$ and $w_n \Rightarrow v$ .

Let $G = (\Phi, \Sigma, R, S)$ be a formal grammar. Then, $L(G) = \{w \in \Sigma^* \mid S \overset{*}{\Rightarrow} w\}$ is the **formal language** generated by G.

An **automaton** is a device that decides, whether a given sentence belongs to a formal language

OSTBAYERISCHE
TECHNISCHE HOCHSCHULE
REGENSBURG

# Generation and Acceptance

Grammar

↓ generates

Language

↓ is accepted by

Automaton

The complexity of the generating grammar influences the complexity of the accepting automaton

# Type-0 Grammar: unrestricted

The mechanism of unrestricted grammars allows the definition of very complex languages that in turn need very complex automata. Restrictions on the form of production rules lead to different types of grammars.

An **unrestricted** formal grammar is called **Type-0 grammar** and can be accepted by a **Turing Machine**. It produces **recursively enumerable** languages.

OTH OSTBAYERISCHE
TECHNISCHE HOCHSCHULE
REGENSBURG

# Type-1 Grammar: context-sensitive

A grammar *G = (Φ,Σ,R,S)* is **context-sensitive**, iff all production rules in R obey the form

- either αAγ→αβγ with α,β,γ ∈ (Φ∪Σ)\*, A∈Φ, β≠ε

- or S→ε.

If S→ε, then S cannot appear in the righthand side of rules in R.

The language of a type-1 grammar is accepted by a **linear bounded automaton** (a nondeterministic Turing machine whose tape is bounded by a constant times the length of the input).

OSTBAYERISCHE
TECHNISCHE HOCHSCHULE
REGENSBURG

# Type−2 Grammar: context−free

A grammar $G = (\Phi, \Sigma, R, S)$ is **context-free**, iff all production rules in $R$ obey the form A$\rightarrow$α with A$\in\Phi$, α $\in (\Phi\cup\Sigma)^*$.

Context-free languages are closed under the following operations:

- Union: if $F$ and $G$ are context-free, so is $F\cup G$

- Concatenation: if $F$ and $G$ are context-free, so is $F\bullet G$

- Kleene Star: if $G$ is context-free, so is G*
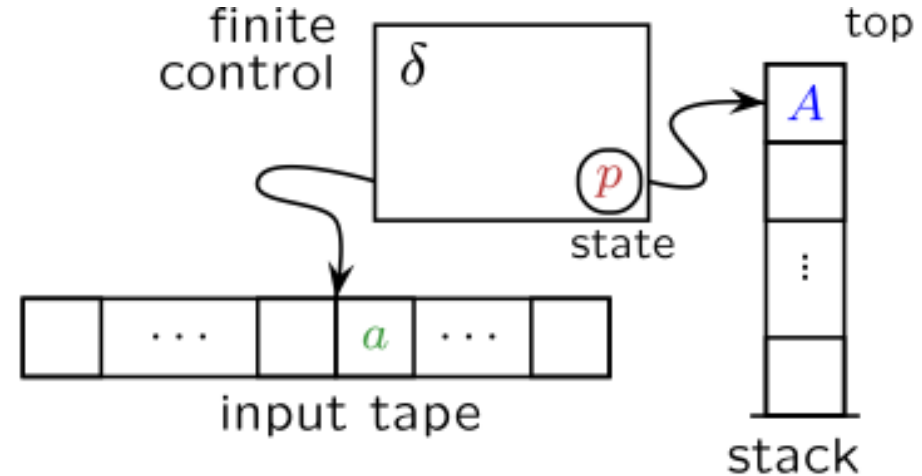
Context-free languages are not closed under the following operations:

- Intersection: from $F$ and $G$ are context-free does not follow that $F\cap G$ is.

- Complement: from $G$ context-free does not follow that $\neg G$ is.

- Difference: from $F$ and $G$ are context-free does not follow that $F\backslash G$ is.

# Pushdown Automaton (=Kellerautomat)

Context-free grammars are accepted by **non-deterministic pushdown automata**.

A non-deterministic push-down automaton PDA=($\Phi,\Sigma,\Delta,\square,\delta,S,F$) consists of

- Alphabet $\Phi$ of states

- Alphabet $\Sigma$ of input symbols, disjunct with $\Phi$

- Alphabet $\Delta$ of stack symbols, disjunct with $\Phi$

- initial stack symbol $\square$

- transition relaton $\delta$: $\Phi\times\Sigma\times(\Delta\cup\square)\rightarrow\rho(\Phi\times\Delta^*)$

- start state $S\in\Phi$

- set of final states $F\subset\Phi$



OSTBAYERISCHE
TECHNISCHE HOCHSCHULE
REGENSBURG

# Type-3 Grammar: regular, left/right linear

A grammar $G = (\Phi, \Sigma, R, S)$ is **right linear (left linear)**, iff all the production rules in $R$ obey the forms

- $A \rightarrow w$

- $A \rightarrow wB$ (left linear: $A \rightarrow Bw$ )

with $A, B \in \Phi$ and $w \in \Sigma^*$.

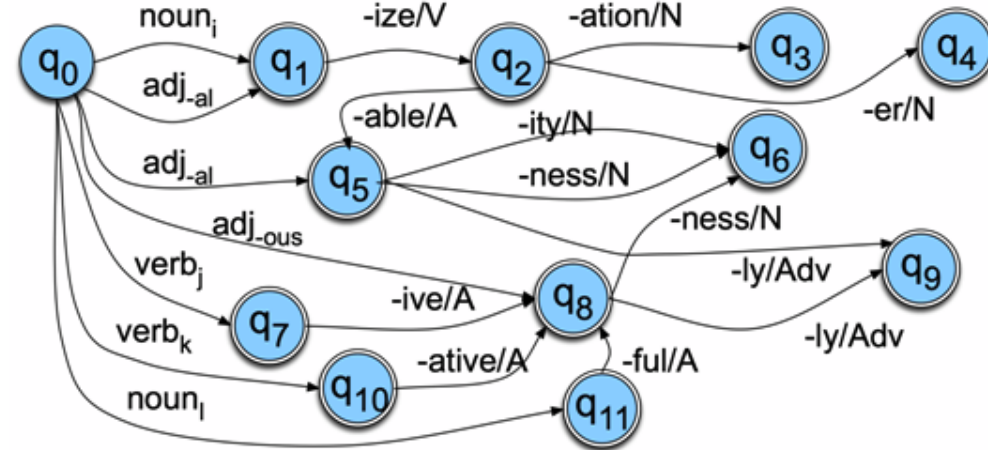Left (right) linear grammars generate **regular** languages:

- $\varnothing$ is regular

- $\{a_i\}$ is regular for $a_i \in$ alphabet $\Sigma$

- if the sets $L_1$ and $L_2$ are regular, then $(L_1 \cap L_2)$ is regular

- if the sets $L_1$ and $L_2$ are regular, then $(L_1 \bullet L_2)$ is regular

- if set L is regular, then also $L^*$.

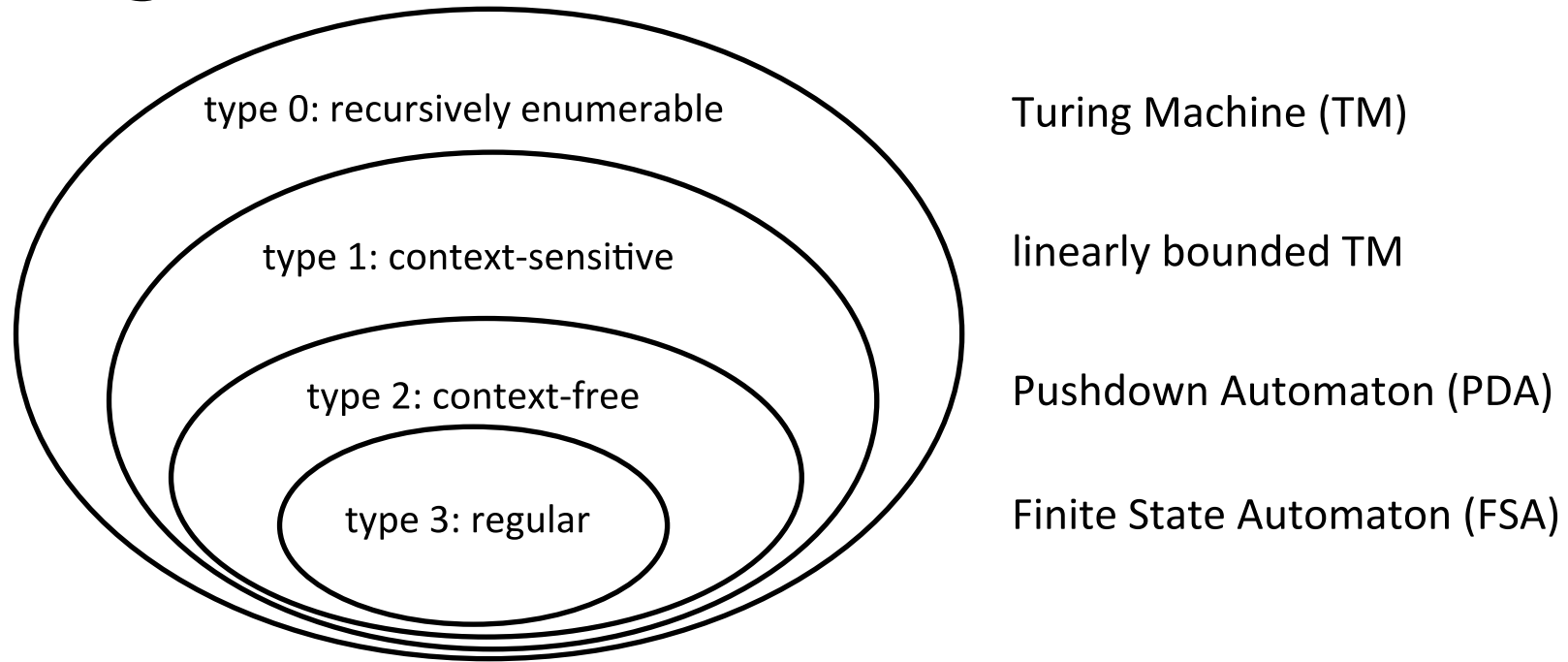These languages can be described through **regular expressions**

... and are accepted by **finite state automata**

OSTBAYERISCHE
TECHNISCHE HOCHSCHULE
REGENSBURG

# Finite State Automaton

- an FSA is a 5-tupel FSA=$(\Phi,\Sigma,\delta,S,F)$ consisting of:
  - a set of states $\Phi$
  - input alphabet $\Sigma$, $\Phi \cap \Sigma = \varnothing$
  - one start state $S \subset \Phi$
  - set of final states $F \subset \Phi$
  - transition function $\delta \subseteq \Phi \times (\Sigma \cup \{\epsilon\}) \rightarrow \Phi$
    - with $\epsilon$ being the "spontaneous transition" that doesn't consume a symbol

# The Chomsky hierarchy of formal languages

type 0: recursively enumerable

Turing Machine (TM)

type 1: context-sensitive

linearly bounded TM

type 2: context-free

Pushdown Automaton (PDA)

type 3: regular

Finite State Automaton (FSA)

- The different classes are proper subsets of each other: the expressivity of type-(n) grammars is truly smaller than type-(n-1) grammars.

- Several other classes are known, e.g. corresponding to deterministic context-free grammars, tree adjoining grammars …

# Principle of Compositionality

- Speech sounds (or letters) form words
- Words follow each other, forming sequences
- sequences of words have an inner structure (syntax)
- meaning (semantics) can be derived from the structure

OSTBAYERISCHE
TECHNISCHE HOCHSCHULE
REGENSBURG

# Means of Compositionality

- differencces across linguistic layers
  - morphology is (mostly) regular
  - syntax uses (mainly context free) grammars
  - semantics can be expressed through first order logics or relation graphs

OSTBAYERISCHE
TECHNISCHE HOCHSCHULE
REGENSBURG

# Principle of Compositionality

- Speech sounds (or letters) form words
  - there are language-specific rules for how speech sounds can be combined (=phonotactics)
  - German allows consonant clusters like in "Strumpf" /ʃ t r ʊ m p f/, where most languages are less liberal.
  - Then again, German sometimes requires "ch" to be /ç/ as in "Nächte" or /x/ as in "Nacht". Try saying this the opposite way, it will feel weird.
  - can you deduce the rule for when "ch" is /ç/ vs. /x/?

# Morphology

- words have internal structure
  - base form

    + morphological properties
  - "largest" → large+Adj+Sup
  - "better" → good+Adj+Comp (exceptions from rules)
  - "runs" → run+Verb+Present+3sg
    OR run+Noun+Pl (ambiguity)

OTH OSTBAYERISCHE TECHNISCHE HOCHSCHULE REGENSBURG

# Morphology

- rules are (by and large) regular, thus, a finite state automaton (FSA) is appropriate

  (wanna crack a nut? → use only a small sledgehammer)

- task for morphologic models
  - given a full word, decompose it into morphemes
  - given base form+morphemic information, generate the resulting "surface structure"

# Why build a model at all?

We could simply store all word forms in a table.

What can be advantages of a morphological model?

# Finite State Transducer

- a finite state **transducer** is a 6-tupel FST=$(\Phi,\Sigma,\Gamma,\delta,S,F)$ consisting of:

    - einer Menge Zustände $\Phi$

    - Eingabealphabet $\Sigma$

    - Ausgabealphabet $\Gamma$

    - Startzuständen $S{\subset}\Phi$

    - Terminalzuständen $F{\subset}\Phi$

    - Übergangsfunktion $\delta{\subseteq}\Phi{\times}(\Sigma{\cup}\{\epsilon\}){\times}(\Gamma{\cup}\{\epsilon\}){\times}\Phi$

        - with $\epsilon$ being the "spontaneous transition"

# FSTs

- FSTs correspond to a finite state automaton but have two alphabets ($\Sigma$ and $\Gamma$), which are used on two different "tapes" (for input and output)

- an FST transduces an input string x to an output string y if there is a sequence of transitions that starts in the start state, ends in a final state and has x as its input and y as its output string

- FSTs can be composed from simpler FSTs just like FSAs (cmp. regular expressions)

# examples of FSTs (for morphology)

# Dealing with Ambiguity through Nondeterminism

- language is ambiguous (on most levels); we embrace nondeterminism as a mechanism to reflect this

- as long as we can't resolve ambiguity, we carry along multiple hypotheses in parallel
  - for FSA: we don't know what path we've taken
  - for FST: different paths produce different output strings

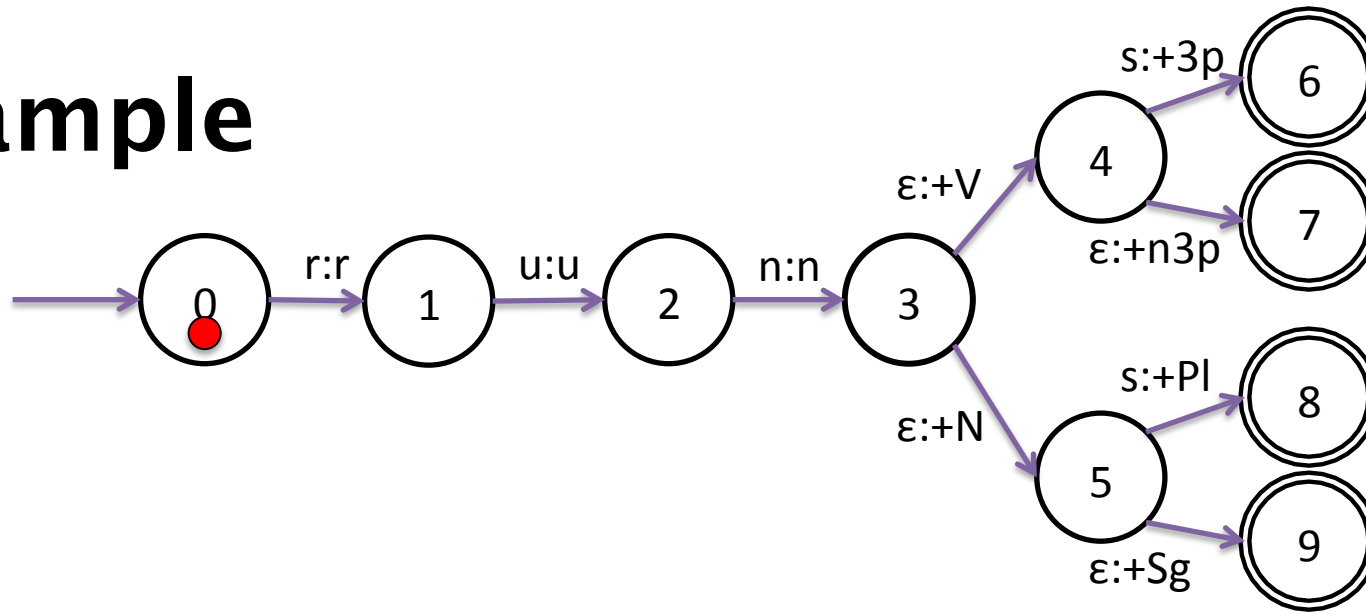- we keep a **set** of possible current states

- acceptance of FST if there's at least one path to a final state

# Example



input
string

| | r | u | n | s | | |
|---|---|---|---|---|---|---|

# Example
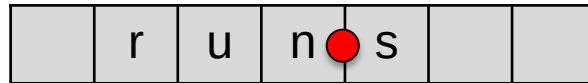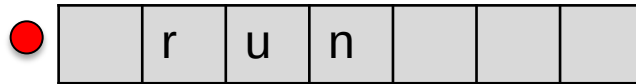


Dots: Keep track of current state and output generated so far.

OSTBAYERISCHE
TECHNISCHE HOCHSCHULE
REGENSBURG

# Example



Transition: dot moves on input tape and to next state, generating output.
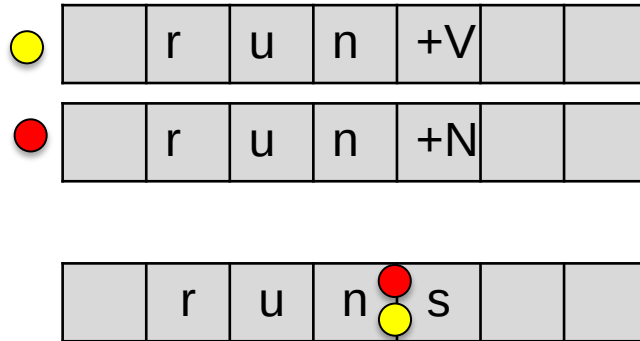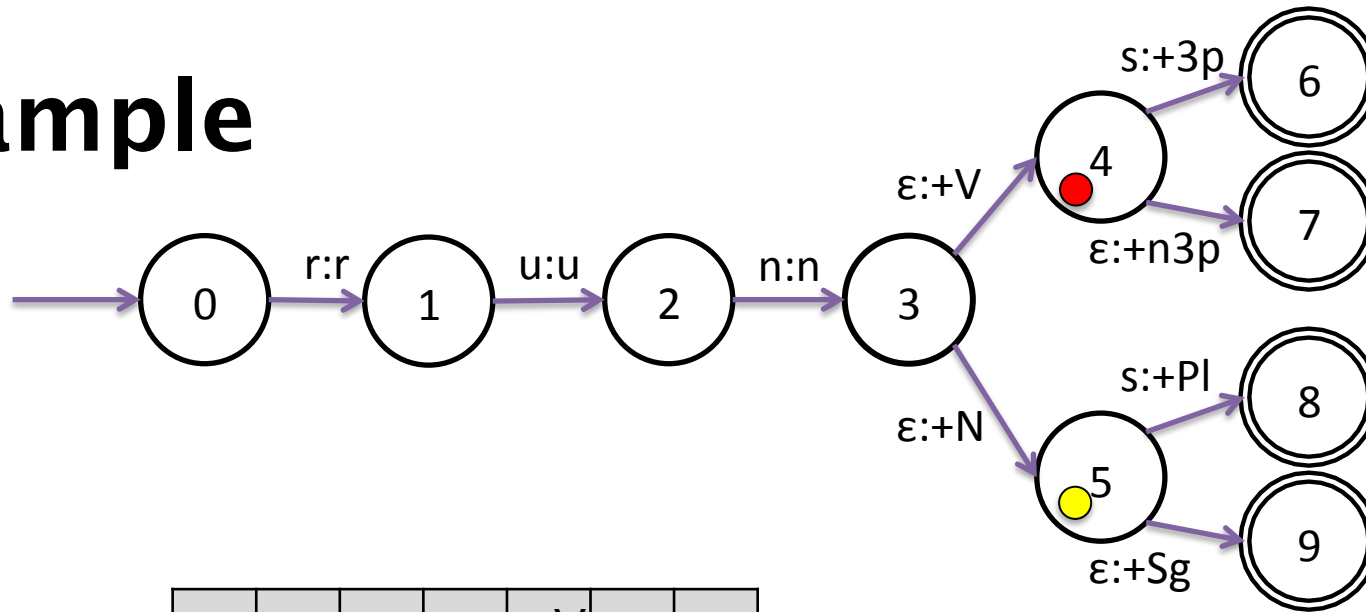
# Example



Transition: dot moves on input tape and to next state, generating output.

OSTBAYERISCHE
TECHNISCHE HOCHSCHULE
REGENSBURG

# Example



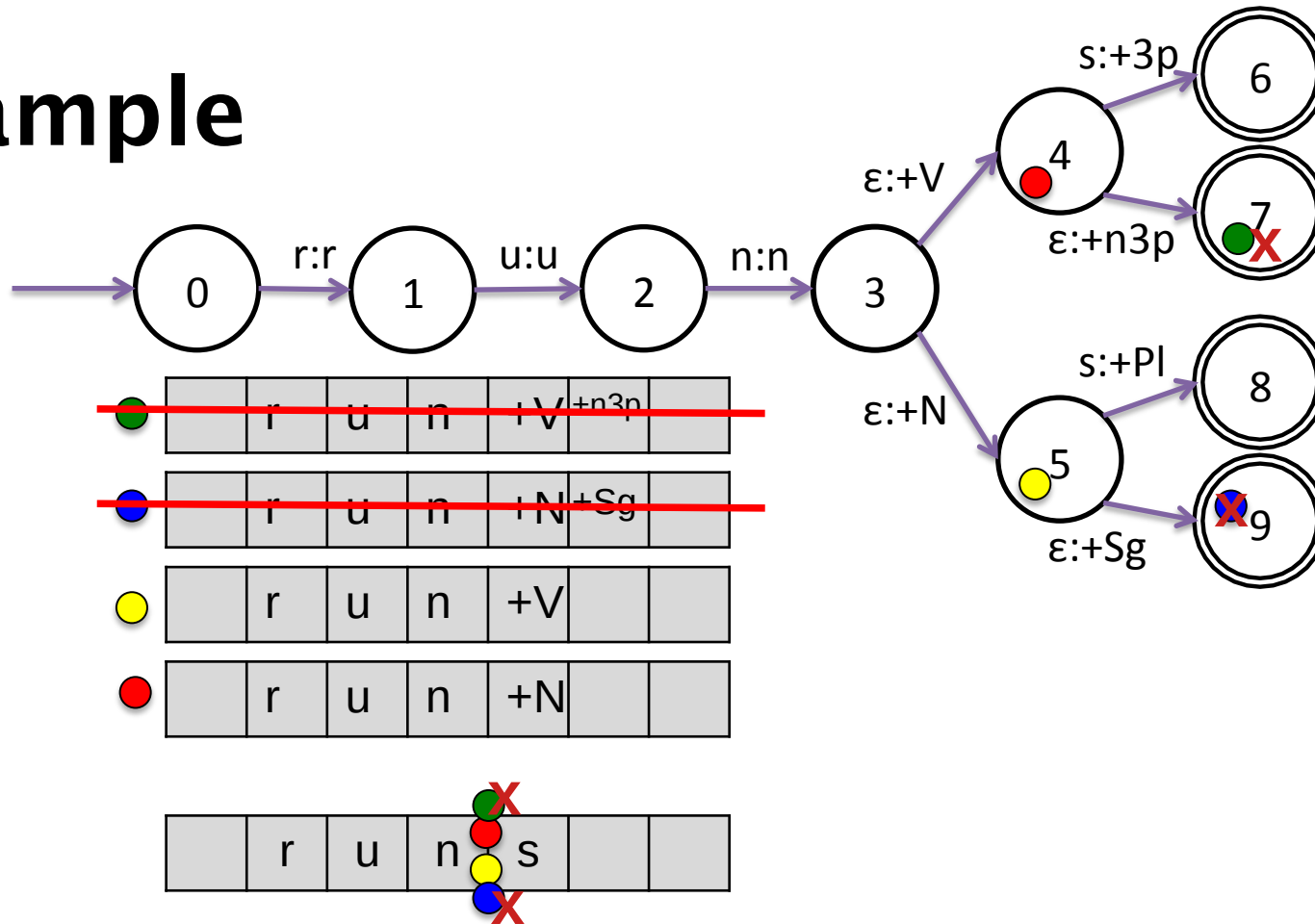Non-determinism: Dot splits. Two possibilities for output tape.

# Example



ε−transitions also introduce nondeterminism.

# Example

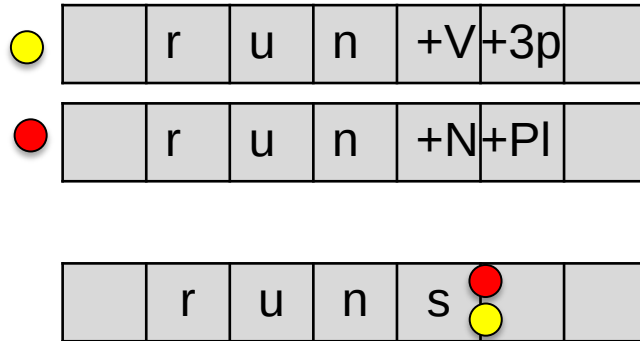

Dots that do not have a follow-up state are removed
as they form no valid hypothesis.

OSTBAYERISCHE
TECHNISCHE HOCHSCHULE
REGENSBURG

# Example



End of input is reached. All dots at final states have successfully transduced.

OSTBAYERISCHE
TECHNISCHE HOCHSCHULE
REGENSBURG

# Summary: FST

- given: some input sequence

- produces: some output sequence
  that is in a regular relation to input sequence

- non-deterministic algorithm to check input/output
  correspondence

  – create all possible outputs given the input

  – typically: weighted FST with beam search (→AGKI)

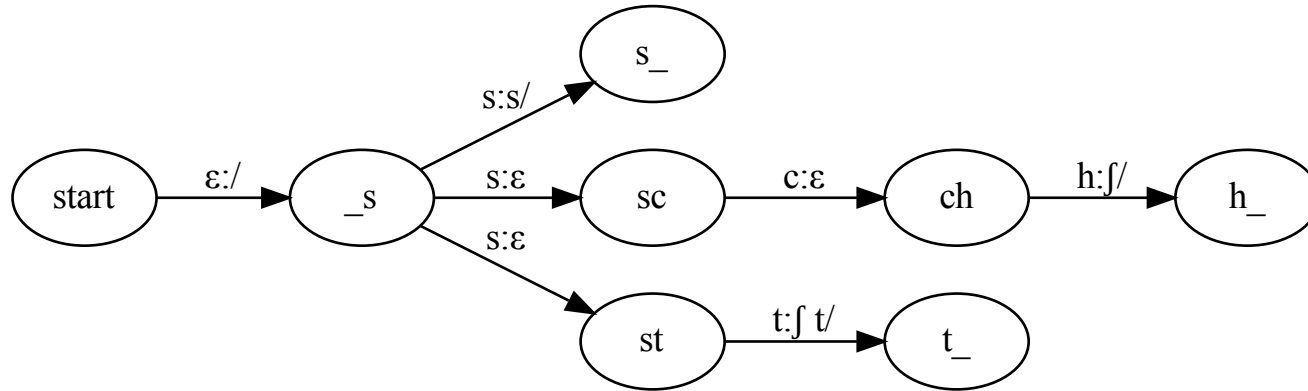- FSTs can be composed & stacked together

# Applications of FSAs/FSTs in NLP

- data structure for lexica (→ spell checking)

- morphological analysis & synthesis

- Segmentation

- Tokenisation

- Sentence boundary detection
  (not every full stop is a sentence boundary)

- Chunk parsing (through cascaded FSTs)

- decoding in speech recognition

- pronunciation modelling (exercise on Friday)

OSTBAYERISCHE
TECHNISCHE HOCHSCHULE
REGENSBURG

# FSTs for spelling

- input alphabet: standard alphabet

- output alphabet: International Phonetic Alphabet, SAMPA, or some other representation for speech sounds

  1. start with some simple rules:
     - "s" → /s/
     - "sch" → /ʃ/
     - "st" → /ʃ t/
     - ...

  2. merge these together to form a complex FST

  - alternative: train from data

OSTBAYERISCHE
TECHNISCHE HOCHSCHULE
REGENSBURG

# FSTs for spelling

# Thank you, your questions?

`timo.baumann@oth-regensburg.de`

# further reading

- Jurafsky and Martin (2009): chapter 3.2
- Carstensen et al. (2004): pp. 198–205

OSTBAYERISCHE
TECHNISCHE HOCHSCHULE
REGENSBURG

# Lehrziele

- die Studierenden kennen automatenbasierte Verfahren in der Sprachverarbeitung

- die Studierenden verstehen den Vorteil von Modellen die sowohl zur Analyse als auch zur Generierung angewendet werden können

- die Studierenden kennen einen Algorithmus zur Dekodierung eines Automaten und können diesen anwenden

- die Studierenden kennen Anwendungen automatenbasierter Verfahren, wie z.B. morphologische Analyse und Aussprachemodellierung

OSTBAYERISCHE
TECHNISCHE HOCHSCHULE
REGENSBURG

41