

Kommunikation mit gesprochener Sprache ...

... ist mehr als Hören und Sprechen:

Kommunikation mit gesprochener Sprache ...

... ist mehr als Hören und Sprechen:

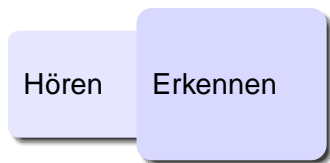


Hören

Kikala brint tovoluti?

Kommunikation mit gesprochener Sprache ...

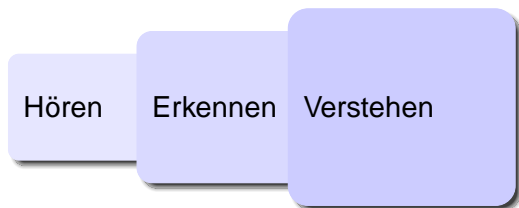
... ist mehr als Hören und Sprechen:



Winter kochtest ganz Blatt?

Kommunikation mit gesprochener Sprache ...

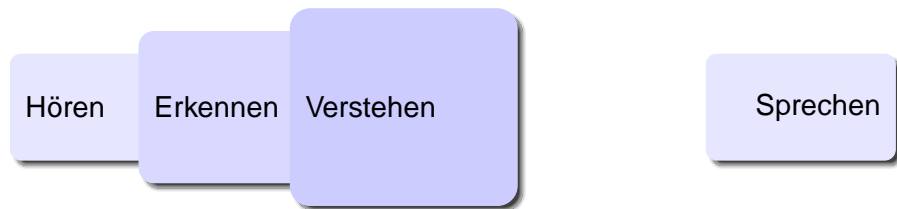
... ist mehr als Hören und Sprechen:



Wann steigt die Party?

Kommunikation mit gesprochener Sprache ...

... ist mehr als Hören und Sprechen:

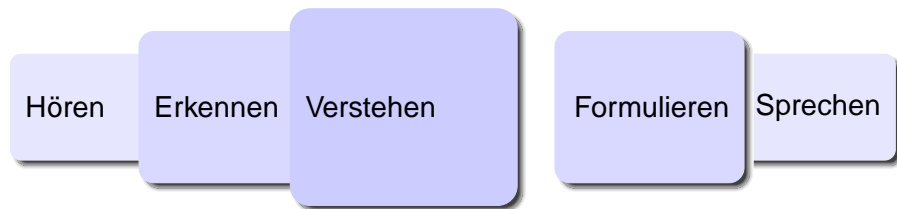


Wann steigt die Party?

Sintu högafi notsi!

Kommunikation mit gesprochener Sprache ...

... ist mehr als Hören und Sprechen:

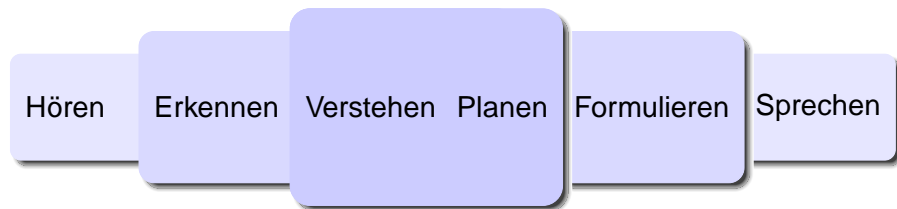


Wann steigt die Party?

Sonderbar werfen die Wellen hinab!

Kommunikation mit gesprochener Sprache ...

... ist mehr als Hören und Sprechen:

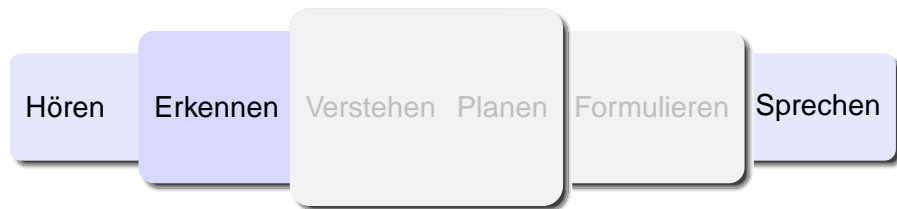


Wann steigt die Party?

Am Freitag, im Hirsch!

Kommunikation mit gesprochener Sprache ...

... ist mehr als Hören und Sprechen:



Wann steigt die Party?

Am Freitag, im Hirsch!

Gesprochene Sprache

Segmentstruktur

Äußerungen Und wie wäre es am Mittwoch?

Gesprochene Sprache

Segmentstruktur

Äußerungen	Und	wie	wäre	es	am	Mittwoch?
		⋮		⋮	⋮	⋮
Phrasen		wie	wäre	es	am	Mittwoch?

Gesprochene Sprache

Segmentstruktur

Äußerungen	Und	wie	wäre	es	am	Mittwoch?
Phrasen		wie	wäre	es	am	Mittwoch?
Wörter		wie	wäre	es		

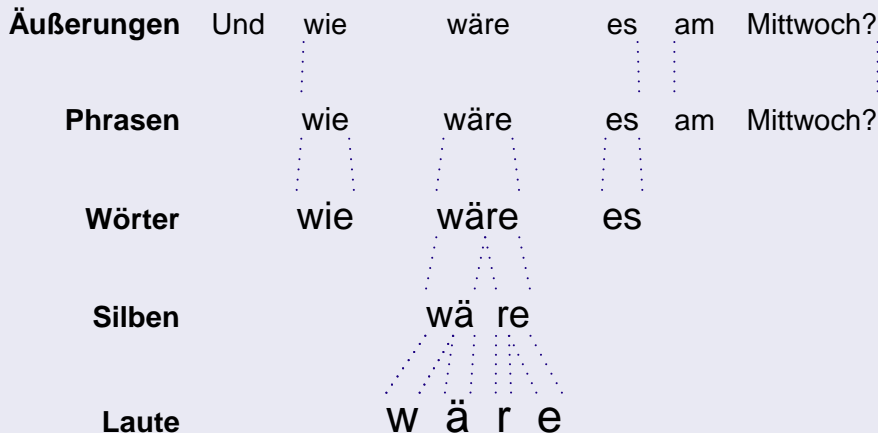
Gesprochene Sprache

Segmentstruktur

Äußerungen	Und	wie	wäre	es	am	Mittwoch?
Phrasen		wie	wäre	es	am	Mittwoch?
Wörter		wie	wäre	es		
Silben			wä re			

Gesprochene Sprache

Segmentstruktur



Gesprochene Sprache

Segmentübergreifende Information (Prosodie)

Gesprochene Sprache

Segmentübergreifende Information (Prosodie)

- Grundfrequenz

Gesprochene Sprache

Segmentübergreifende Information (Prosodie)

- Grundfrequenz
- Rhythmus

Gesprochene Sprache

Segmentübergreifende Information (Prosodie)

- Grundfrequenz
- Rhythmus
- Lautstärke

Gesprochene Sprache

Segmentübergreifende Information (Prosodie)

- Grundfrequenz
- Rhythmus
- Lautstärke

wird verwendet ...

Segmentübergreifende Information (Prosodie)

- Grundfrequenz
- Rhythmus
- Lautstärke

wird verwendet ...

- ... für Hervorhebungen (Neues, Wichtiges, Unerwartetes, ...)

Segmentübergreifende Information (Prosodie)

- Grundfrequenz
- Rhythmus
- Lautstärke

wird verwendet ...

- ... für Hervorhebungen (Neues, Wichtiges, Unerwartetes, ...)
- ... zur Gliederung (Phrasen, Sätze)

Segmentübergreifende Information (Prosodie)

- Grundfrequenz
- Rhythmus
- Lautstärke

wird verwendet ...

- ... für Hervorhebungen (Neues, Wichtiges, Unerwartetes, ...)
- ... zur Gliederung (Phrasen, Sätze)
- ... zur Moduskennzeichnung (Aussage, Frage, ...)

Segmentübergreifende Information (Prosodie)

- Grundfrequenz
- Rhythmus
- Lautstärke

wird verwendet ...

- ... für Hervorhebungen (Neues, Wichtiges, Unerwartetes, ...)
- ... zur Gliederung (Phrasen, Sätze)
- ... zur Moduskennzeichnung (Aussage, Frage, ...)
- ... zur Dialogsteuerung (Vergabe der Initiative)

Segmentübergreifende Information (Prosodie)

- Grundfrequenz
- Rhythmus
- Lautstärke

wird verwendet ...

- ... für Hervorhebungen (Neues, Wichtiges, Unerwartetes, ...)
- ... zur Gliederung (Phrasen, Sätze)
- ... zur Moduskennzeichnung (Aussage, Frage, ...)
- ... zur Dialogsteuerung (Vergabe der Initiative)
- ... zum Ausdruck von Emotionen (Freude, Angst, Überraschung, Verlegenheit, ...)

Sprachsynthese

Zwei Herangehensweisen:

Sprachsynthese

Zwei Herangehensweisen:

Vollsynthese

- Erzeugen des Sprachsignals durch Ton- und Rauschgeneratoren
- akzeptable Verständlichkeit
- geringe Natürlichkeit

Sprachsynthese

Zwei Herangehensweisen:

Vollsynthese

- Erzeugen des Sprachsignals durch Ton- und Rauschgeneratoren
- akzeptable Verständlichkeit
- geringe Natürlichkeit

reproduktive Synthese

- Aufnehmen und Wiedergeben menschlicher Sprachsignale
- hohe Verständlichkeit
- gute bis hohe Natürlichkeit

reproduktive Synthese

Was sind geeignete Basiseinheiten?

- ganze Phrasen: nur für Spezialanwendungen
- Laute: schlechte Qualität
- Kompromiss: flexible Ermittlung der Basiselemente aus Korpusdaten

reproduktive Synthese

Was sind geeignete Basiseinheiten?

- ganze Phrasen: nur für Spezialanwendungen
- Laute: schlechte Qualität
- Kompromiss: flexible Ermittlung der Basiselemente aus Korpusdaten

Wie werden die Basiselemente verkettet?

- harter Schnitt erzeugt Knackgeräusche
- "weiche" Übergänge erforderlich
- prosodische Variation durch spezielle Transformationsverfahren

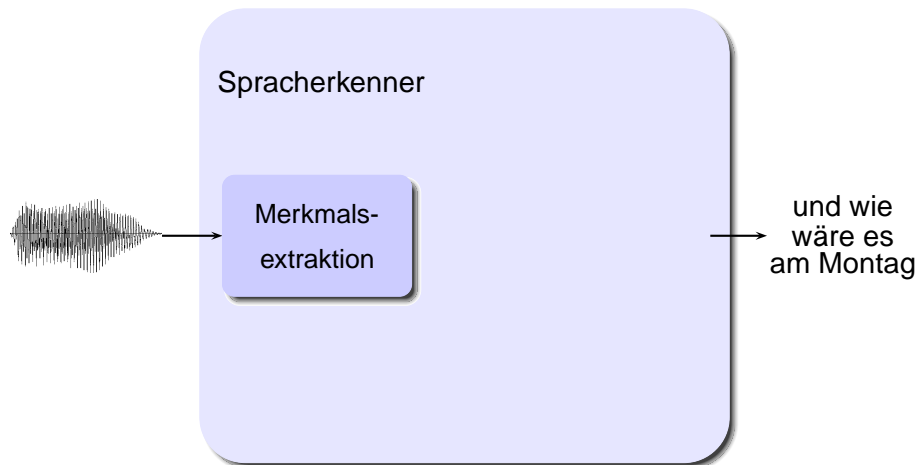
Spracherkennung

- nur Berücksichtigung von Lautcharakteristika
- "Training" von Modellen auf großen Sprachdatensammlungen
- Vernachlässigung der Prosodie
- nur Erkennung, kein Sprachverstehen!

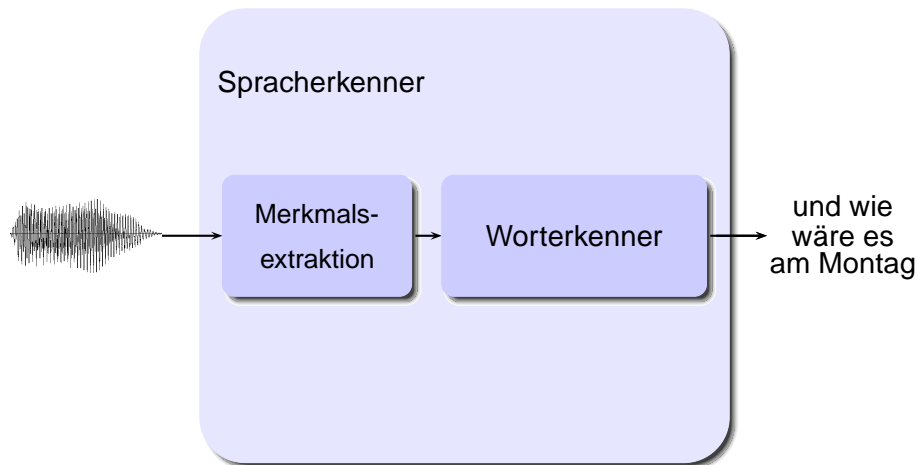
Spracherkennung



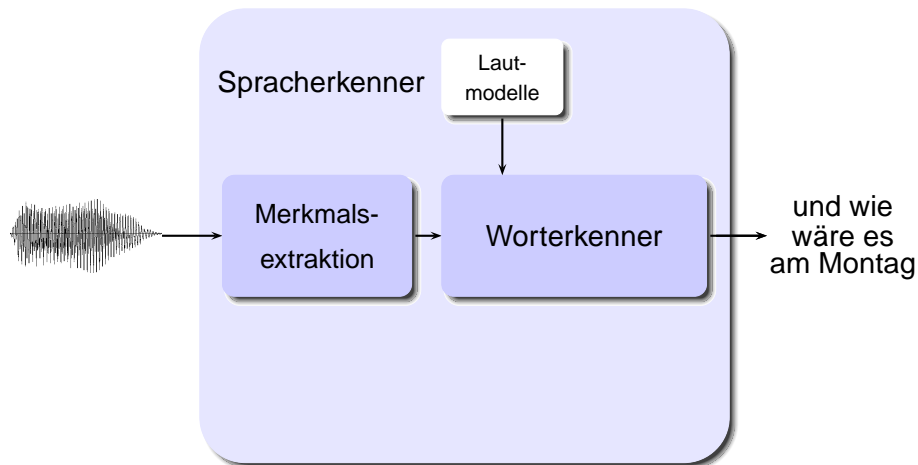
Spracherkennung



Spracherkennung



Spracherkennung



Spracherkennung

Spracherkennung

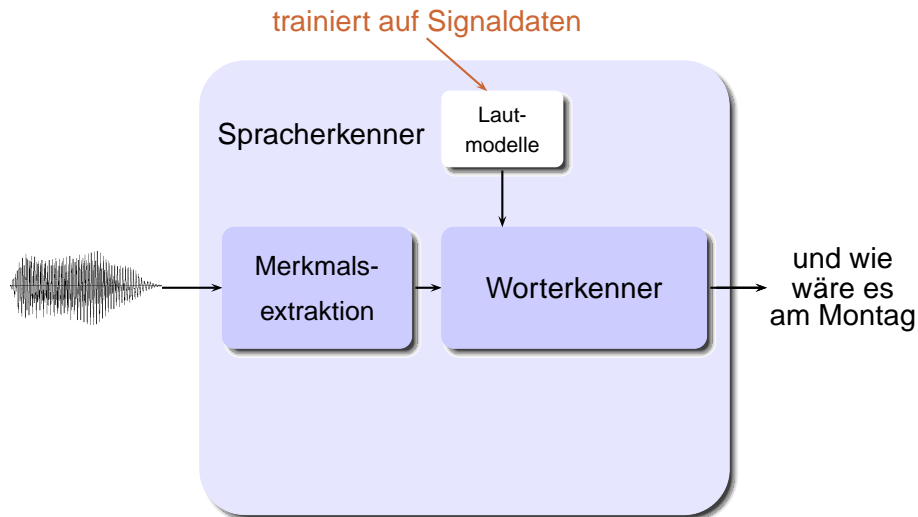
Laut-
modelle

- Modelle für jeden Laut im Kontext seiner Nachbarlaute
 $m-a+m$, $m-a+n$, $d-a+n$, ...
- Berechnung der Wahrscheinlichkeit, dass das Sprachsignal durch das Modell erzeugt wurde
- Zustände, Zustandsübergänge
- Transitionswahrscheinlichkeiten
- Emissionswahrscheinlichkeiten

erkenner

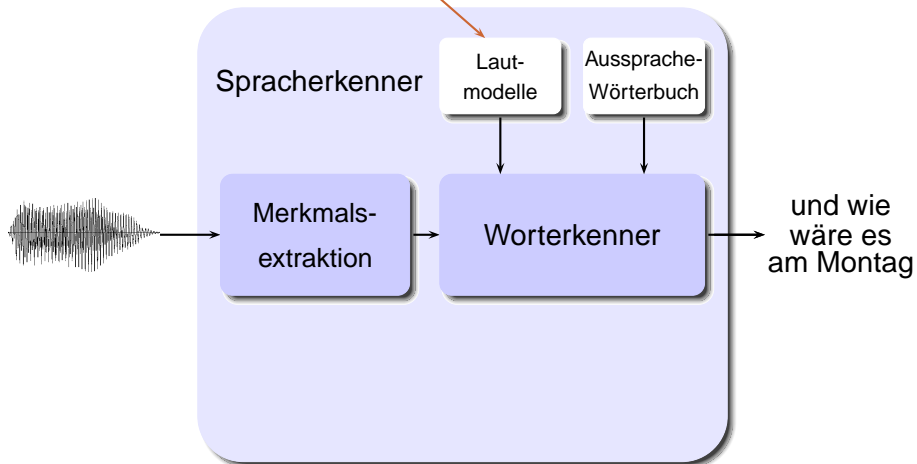
und wie
wäre es
am Montag

Spracherkennung



Spracherkennung

trainiert auf Signaldaten



Spracherkennung

trainiert auf Signaldaten

- eine oder mehrere Lautfolgen für jede Wortform

Mittwoch m i t v o x sp
wäre v e h r 2 sp

- Verkettung von Lautmodellen zu Wortmodellen

Mittwoch:

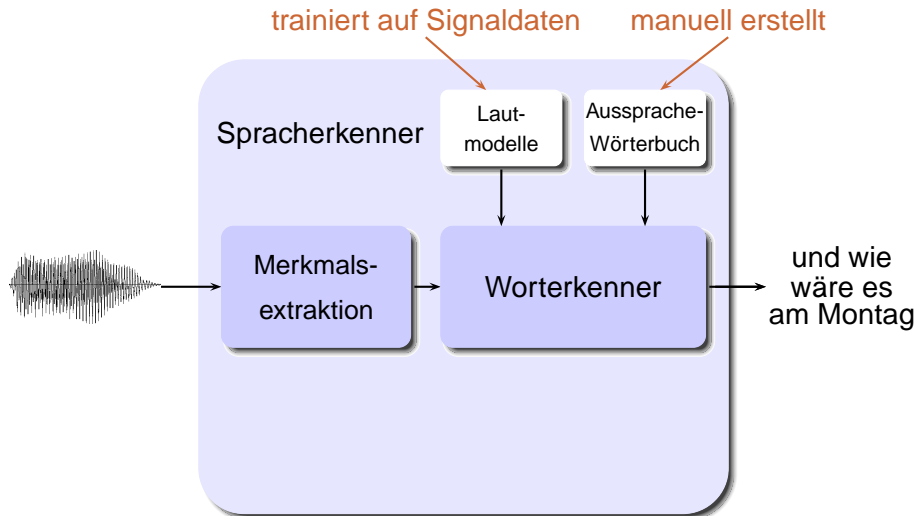
sp-m+i m-i+t i-t+v t-v+o ...

Aussprache-
Wörterbuch

Wörterkennner

und wie
wäre es
am Montag

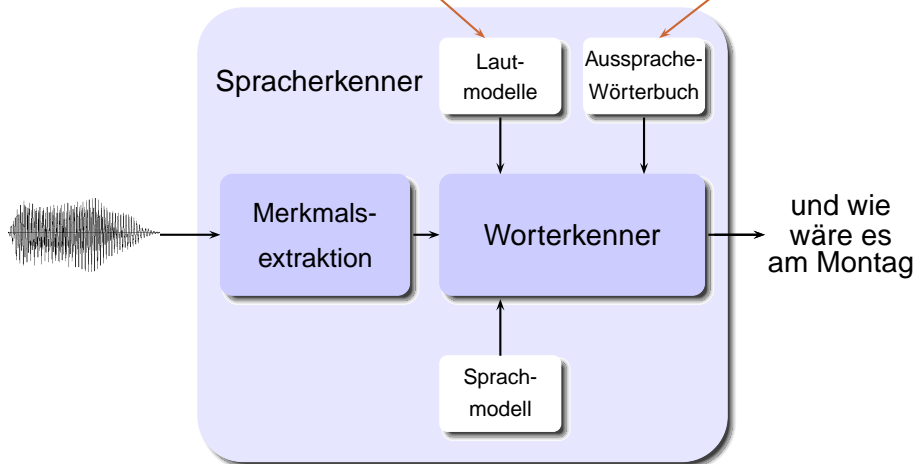
Spracherkennung



Spracherkennung

trainiert auf Signaldaten

manuell erstellt



Spracherkennung

trainiert auf Signaldaten

manuell erstellt

- Berechnung der Wahrscheinlichkeit für komplette Äußerungen
- Wahrscheinlichkeiten für Wortpaare, -tripel oder -quadrupel
 - $p(\text{wir}|\text{dann wollen})$
 - $p(\text{Mittwoch}|\text{dann wollen})$
- wenig geeignet für Dialogsysteme

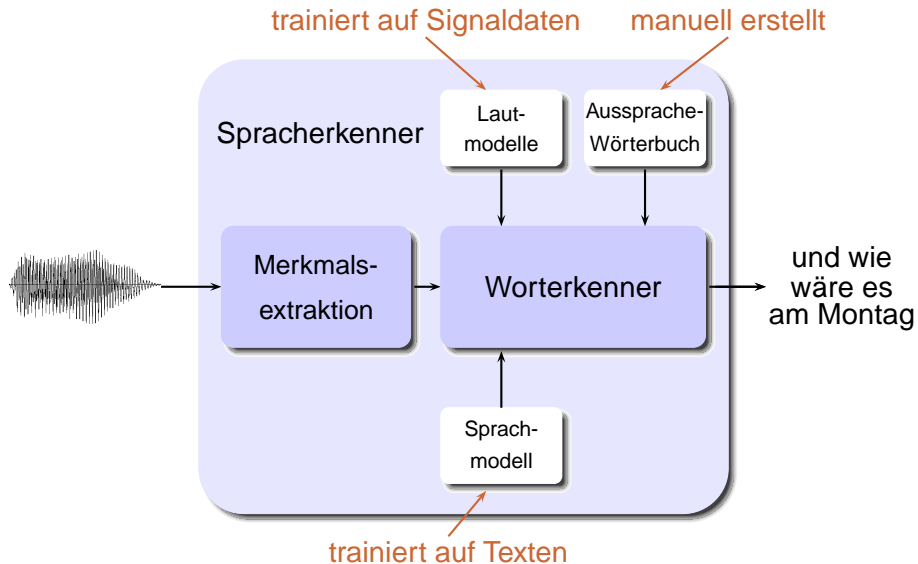
Aussprache-
Wörterbuch

erkenner

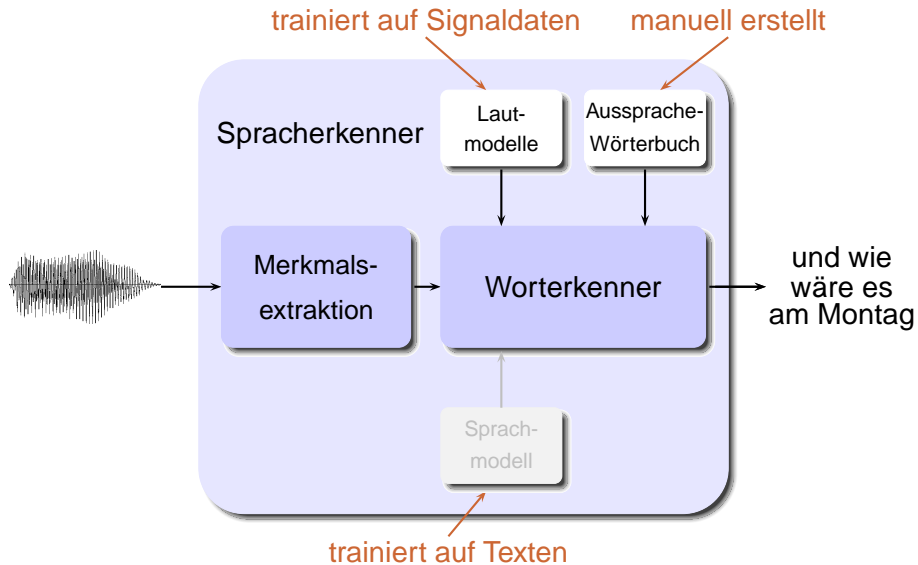
und wie
wäre es
am Montag

Sprach-
modell

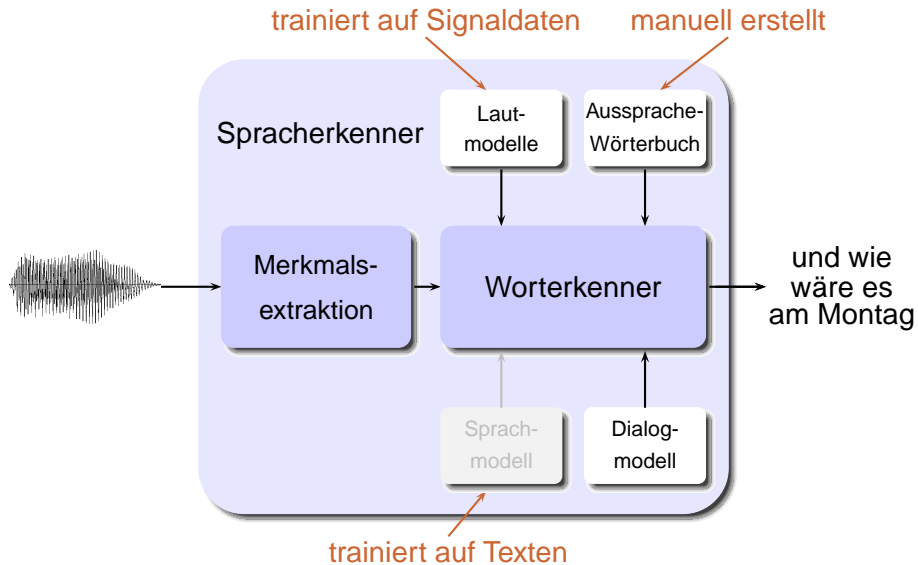
Spracherkennung



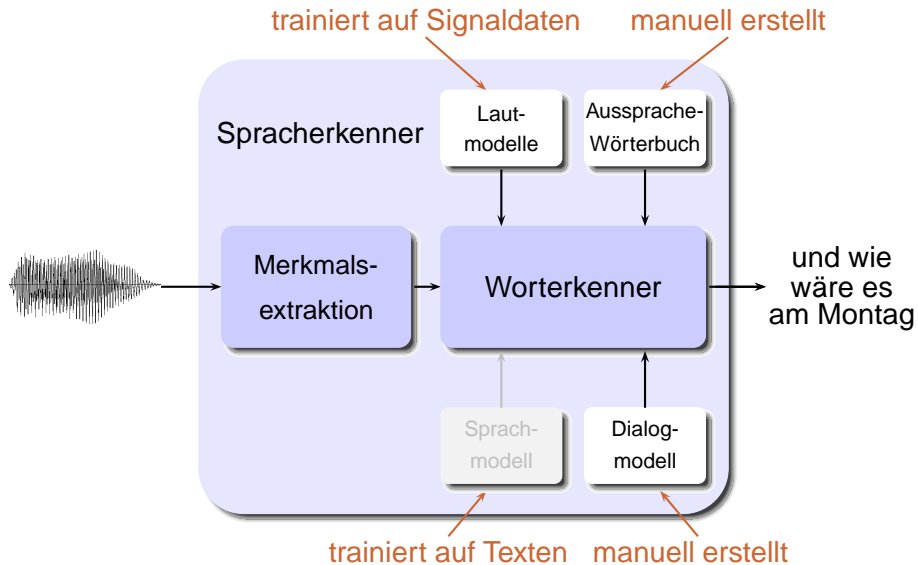
Spracherkennung



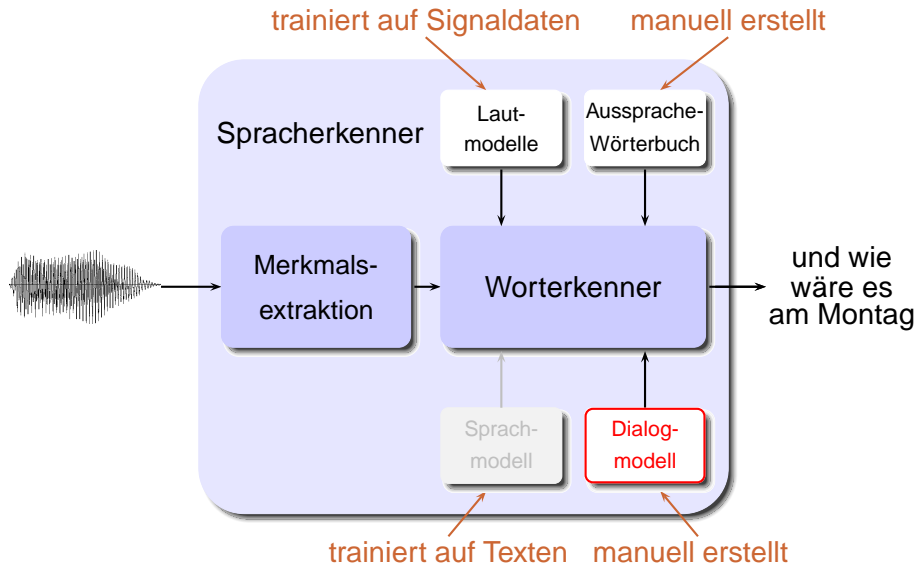
Spracherkennung



Spracherkennung



Spracherkennung



Dialogmodellierung

- dynamische Einschränkung des Erkennerwortschatzes in Abhängigkeit vom Dialogzustand

Dialogmodellierung

- dynamische Einschränkung des Erkennerwortschatzes in Abhängigkeit vom Dialogzustand
- Wozu braucht man das?

Dialogmodellierung

- dynamische Einschränkung des Erkennerwortschatzes in Abhängigkeit vom Dialogzustand
- Wozu braucht man das?
 - Erkennungssicherheit erhöhen
→ Was wurde gesagt?

Dialogmodellierung

- dynamische Einschränkung des Erkennerswortschatzes in Abhängigkeit vom Dialogzustand
- Wozu braucht man das?
 - Erkennungssicherheit erhöhen
 - Was wurde gesagt?
 - ähnliche Aussprache: *Mai* oder *drei*?
 - verschiedene Sprecher
 - schlechte Übertragungsqualität

Dialogmodellierung

- dynamische Einschränkung des Erkennerswortschatzes in Abhängigkeit vom Dialogzustand
- Wozu braucht man das?
 - Erkennungssicherheit erhöhen
 - Was wurde gesagt?
 - ähnliche Aussprache: *Mai* oder *drei*?
 - verschiedene Sprecher
 - schlechte Übertragungsqualität
 - semantische Interpretation erleichtern
 - Was wird von der Maschine erwartet?

Dialogmodellierung

- dynamische Einschränkung des Erkennerwortschatzes in Abhängigkeit vom Dialogzustand
- Wozu braucht man das?
 - Erkennungssicherheit erhöhen
 - Was wurde gesagt?
 - ähnliche Aussprache: *Mai* oder *drei*?
 - verschiedene Sprecher
 - schlechte Übertragungsqualität
 - semantische Interpretation erleichtern
 - Was wird von der Maschine erwartet?
 - Zahl → Geldbetrag, Uhrzeit, Datum, Kontonummer, ...
 - Mehrdeutigkeit: *Wann werden S/sie kommen?*
 - Referenzauflösung: Was wird durch *sie* bezeichnet?

Dialogmodellierung

- dynamische Einschränkung des Erkennerwortschatzes in Abhängigkeit vom Dialogzustand
- Wozu braucht man das?
 - Erkennungssicherheit erhöhen
 - Was wurde gesagt?
 - ähnliche Aussprache: *Mai* oder *drei*?
 - verschiedene Sprecher
 - schlechte Übertragungsqualität
 - semantische Interpretation erleichtern
 - Was wird von der Maschine erwartet?
 - Zahl → Geldbetrag, Uhrzeit, Datum, Kontonummer, ...
 - Mehrdeutigkeit: *Wann werden S/sie kommen?*
 - Referenzauflösung: Was wird durch *sie* bezeichnet?
- eine Maschine hat keinen gesunden Menschenverstand!

Dialogmodellierung

- Dialogzustände: Aufforderung zur Eingabe (Prompt)
- Übergänge zwischen Dialogzuständen: Erkennung von Nutzeräußerungen

Dialogmodellierung

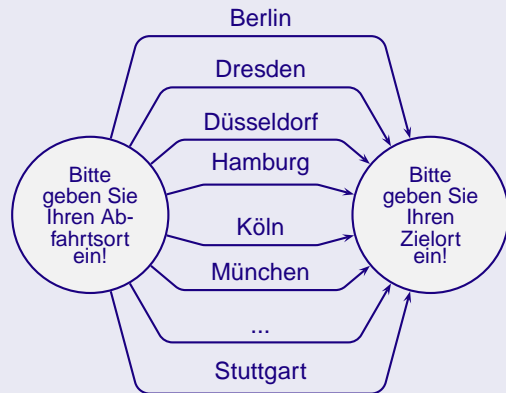
- Dialogzustände: Aufforderung zur Eingabe (Prompt)
- Übergänge zwischen Dialogzuständen: Erkennung von Nutzeräußerungen



Bitte
geben Sie
Ihren Ab-
fahrtsort
ein!

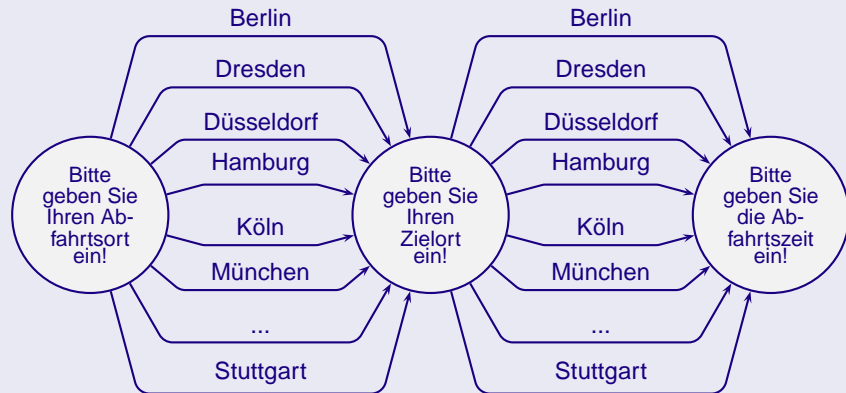
Dialogmodellierung

- Dialogzustände: Aufforderung zur Eingabe (Prompt)
- Übergänge zwischen Dialogzuständen: Erkennung von Nutzeräußerungen



Dialogmodellierung

- Dialogzustände: Aufforderung zur Eingabe (Prompt)
- Übergänge zwischen Dialogzuständen: Erkennung von Nutzeräußerungen



Dialogmodellierung

- Mehrfachverwendung von Teilnetzen



Bitte
geben Sie
Ihren Ab-
fahrtsort
ein!

Dialogmodellierung

- Mehrfachverwendung von Teilnetzen



Dialogmodellierung

- Mehrfachverwendung von Teilnetzen



Dialogmodellierung

- sprecherunabhängige Spracherkennung ist unsicher
- insbesondere bei Telefoneingabe
- Erhöhen der Zuverlässigkeit durch Rückfragen



Bitte
geben Sie
Ihren Ab-
fahrtsort
ein!

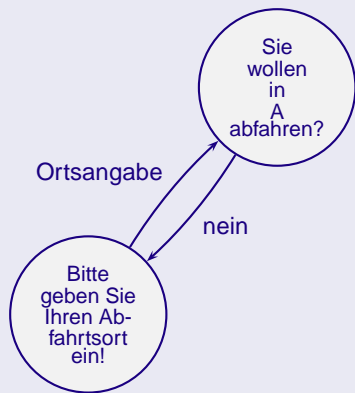
Dialogmodellierung

- sprecherunabhängige Spracherkennung ist unsicher
- insbesondere bei Telefoneingabe
- Erhöhen der Zuverlässigkeit durch Rückfragen



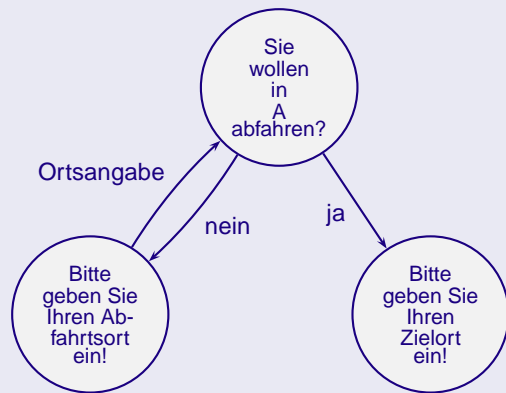
Dialogmodellierung

- sprecherunabhängige Spracherkennung ist unsicher
- insbesondere bei Telefoneingabe
- Erhöhen der Zuverlässigkeit durch Rückfragen



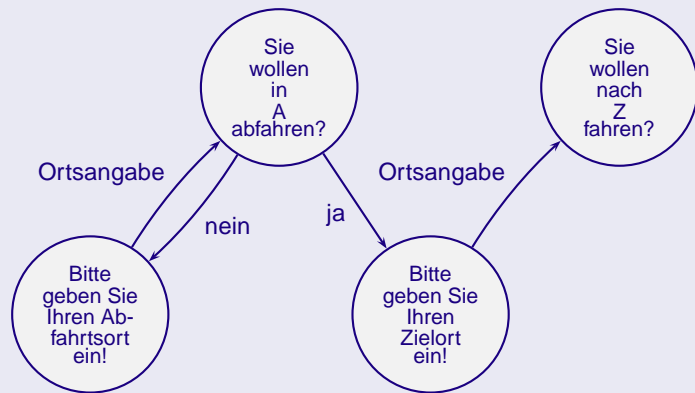
Dialogmodellierung

- sprecherunabhängige Spracherkennung ist unsicher
- insbesondere bei Telefoneingabe
- Erhöhen der Zuverlässigkeit durch Rückfragen



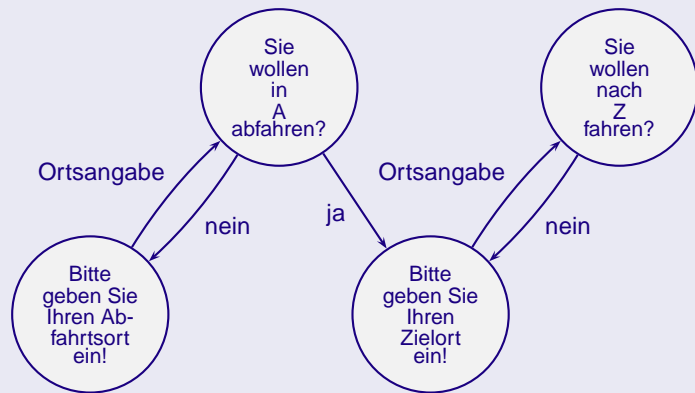
Dialogmodellierung

- sprecherunabhängige Spracherkennung ist unsicher
- insbesondere bei Telefoneingabe
- Erhöhen der Zuverlässigkeit durch Rückfragen



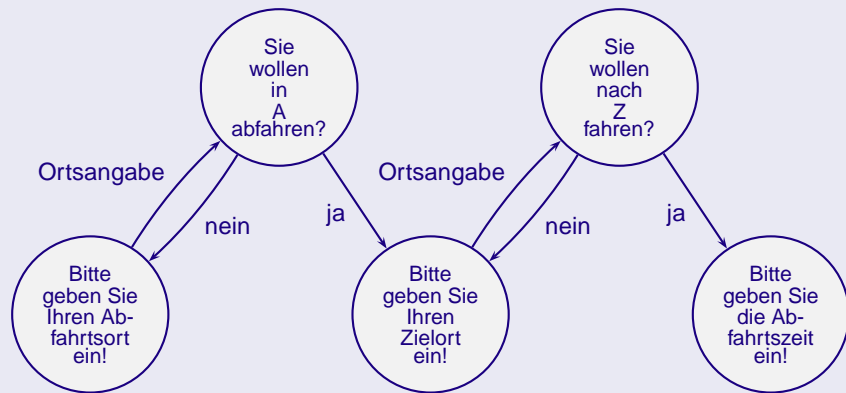
Dialogmodellierung

- sprecherunabhängige Spracherkennung ist unsicher
- insbesondere bei Telefoneingabe
- Erhöhen der Zuverlässigkeit durch Rückfragen



Dialogmodellierung

- sprecherunabhängige Spracherkennung ist unsicher
- insbesondere bei Telefoneingabe
- Erhöhen der Zuverlässigkeit durch Rückfragen



Dialogmodellierung

- theoretische Grundlage: deterministischer endlicher Automat

Dialogmodellierung

- theoretische Grundlage: deterministischer endlicher Automat
 - einfachstes Automatenmodell der Informatik
 - effiziente Implementierung
 - gute Vorhersagefähigkeit → starke Einschränkung des aktiven Wortschatzes

Dialogmodellierung

- theoretische Grundlage: deterministischer endlicher Automat
 - einfachstes Automatenmodell der Informatik
 - effiziente Implementierung
 - gute Vorhersagefähigkeit → starke Einschränkung des aktiven Wortschatzes
- für natürliche Dialogführung zu rigide → Erweiterungen nötig

Dialogmodellierung

- theoretische Grundlage: deterministischer endlicher Automat
 - einfachstes Automatenmodell der Informatik
 - effiziente Implementierung
 - gute Vorhersagefähigkeit → starke Einschränkung des aktiven Wortschatzes
- für natürliche Dialogführung zu rigide → Erweiterungen nötig
 - wechselnde Prompts
 - "Hineinreden" in den Prompt (barge in)
 - ...

Dialogmodellierung

- Dialogmodellierung erfordert Vorhersehen möglicher Nutzerreaktionen
 - Hineinversetzen in den Nutzer
 - Wizard-of-Oz-Experimente

Dialogmodellierung

- Dialogmodellierung erfordert Vorhersehen möglicher Nutzerreaktionen
 - Hineinversetzen in den Nutzer
 - Wizard-of-Oz-Experimente
- Dialogmodellierung schränkt die sprachliche Möglichkeiten eines Nutzers stark ein
 - Lenkung des Nutzers durch Vorgabe zulässiger Äußerungen (akustisches Menü)