

# Dialogue Systems

NPFL123 Dialogové systémy

## 4. Smart Assistants & Question Answering + a little machine learning recap

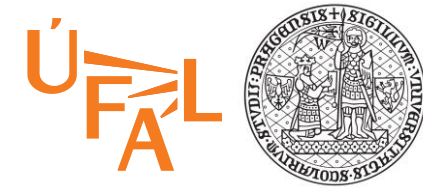
**Ondřej Dušek** & Ondřej Plátek & Jan Cuřín

[ufal.cz/npfl123](http://ufal.cz/npfl123)

12. 3. 2019

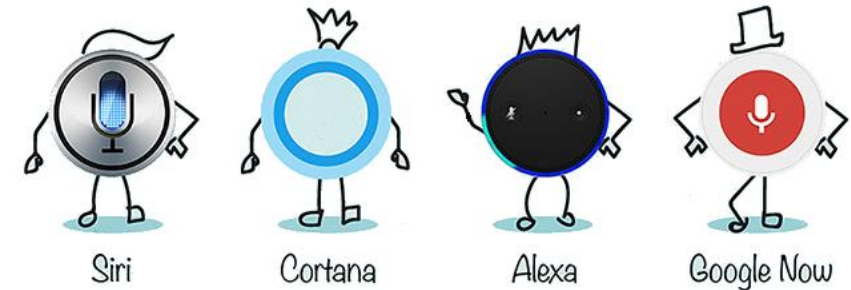
# Virtual Assistants

(voice/smart/conversational assistants)



“Definition”: voice-operated **software** (dialogue system) capable of **answering questions, performing tasks** & basic dialogue in **multiple domains**

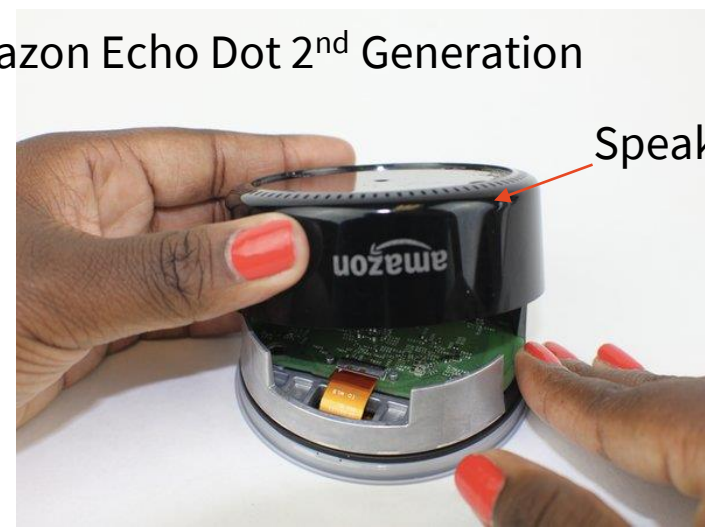
- Apple Siri (2011) –question answering & iOS functions
- Now every major IT company has them
  - Microsoft Cortana (2014)
  - Amazon Alexa (2014)
  - Google Assistant (2016)
  - Samsung Bixby (2017)
  - Mycroft (open-source, 2018)
  - Clova (Naver, 2017) – Korean & Japanese
  - Alice (Yandex, 2017) – Russian
  - DuerOS (Baidu, 2017), AliGenie (Alibaba, 2017) – Chinese



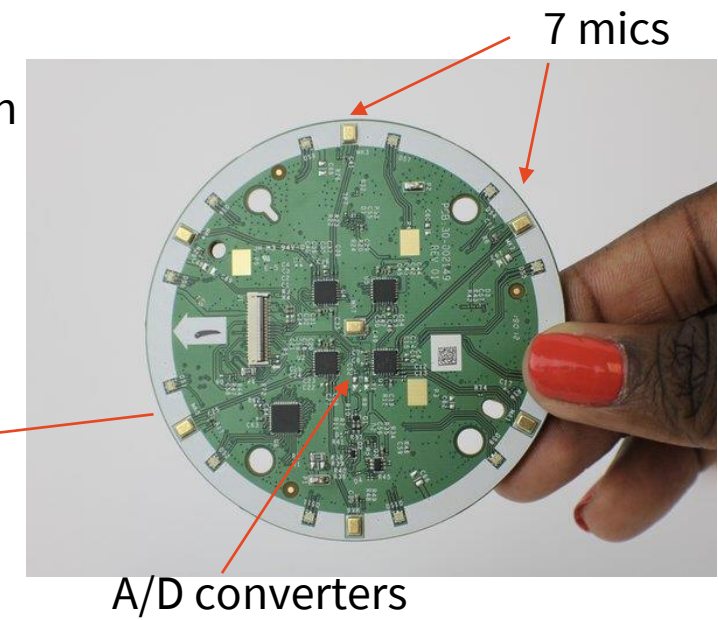
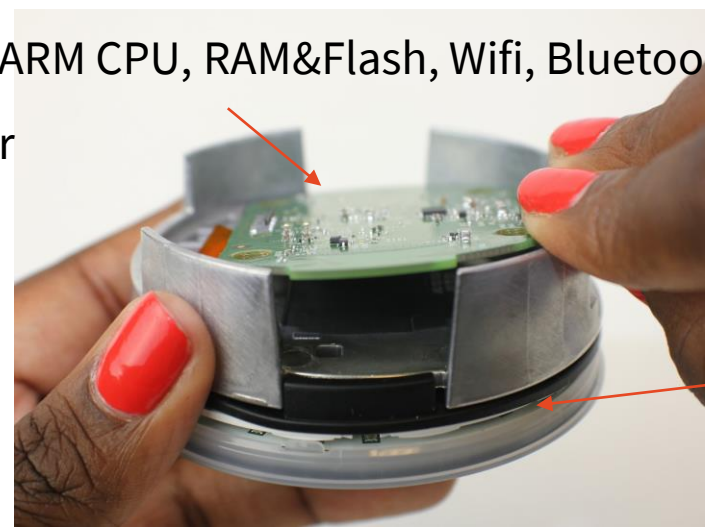
# Smart Speakers

- Internet-connected mic & speaker with a virtual assistant running
  - optionally video (display/camera)
  - ~ same functionality as virtual assistants in phones/computers
  - Amazon Echo (Alexa), Google Home (Assistant), Apple HomePod (Siri) [...]
- Main point: multiple microphones – far-field ASR

Amazon Echo Dot 2<sup>nd</sup> Generation



ARM CPU, RAM&Flash, Wifi, Bluetooth



[https://www.ifixit.com/Device/Amazon Echo Dot 2nd Generation](https://www.ifixit.com/Device/Amazon_Echo_Dot_2nd_Generation)

# Capabilities

- Out of the box:
  - Question answering
  - Web search
  - News & Weather
  - Scheduling
  - Navigation
  - Local information
  - Shopping
  - Media playback
  - Home automation
- a lot of it through 3<sup>rd</sup> party APIs
- the domains are well connected



<https://www.lifehacker.com.au/2018/02/specs-showdown-google-home-vs-amazon-echo-vs-apple-homepod/>

# Demos

Raven H (powered by DuerOS, Baidu)

<https://www.youtube.com/watch?v=iqMjTNjFIMk>



Google Assistant

<https://www.youtube.com/watch?v=JONGt32mfRY>

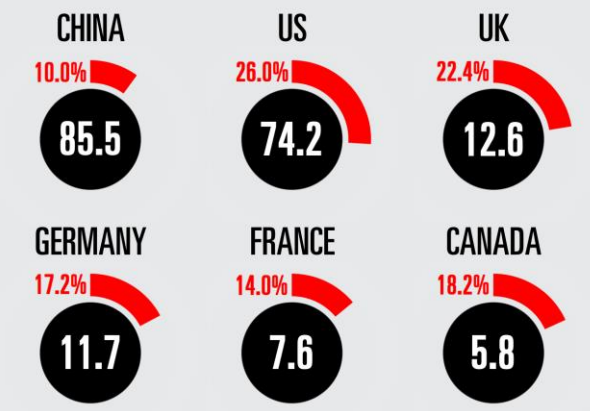




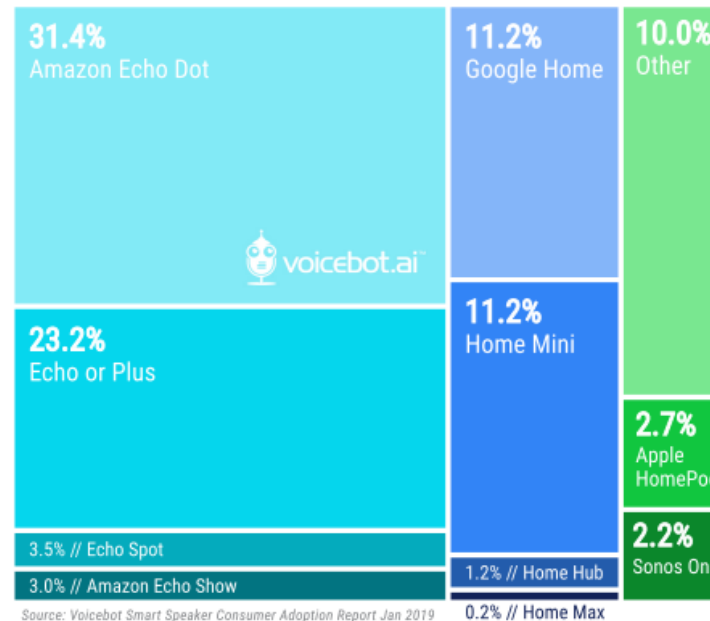
# Smart Speaker Adoption

- >26% US adults have a smart speaker
  - 40% yearly growth in 2018
  - this is very different across the globe
- Amazon leads in the US, Google on the rise

Smart Speaker Users in Select Markets, 2019  
millions and % of internet users



U.S. Smart Speaker Market Share by Device - January 2019

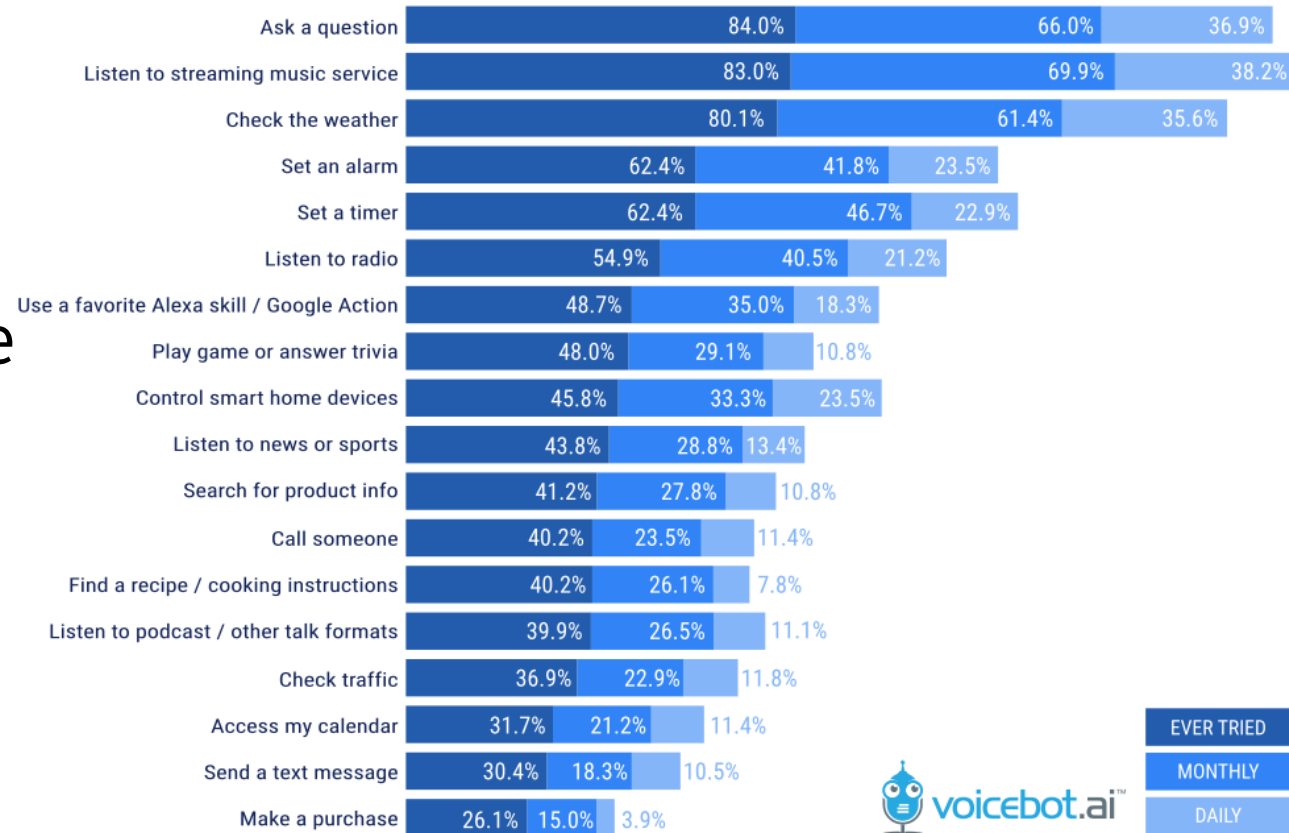


Source: Voicebot Smart Speaker Consumer Adoption Report Jan 2019

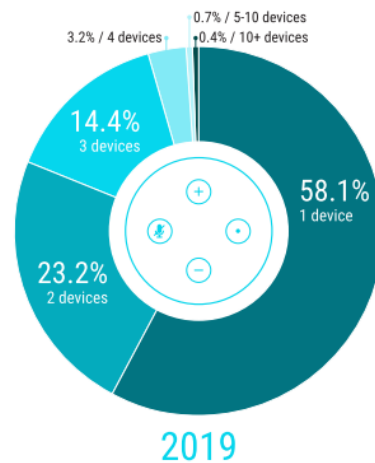
# Smart Speaker Adoption

- People really use them
  - early adopters – more intensively
  - correlated with phone assistant usage
- Many people have more than one

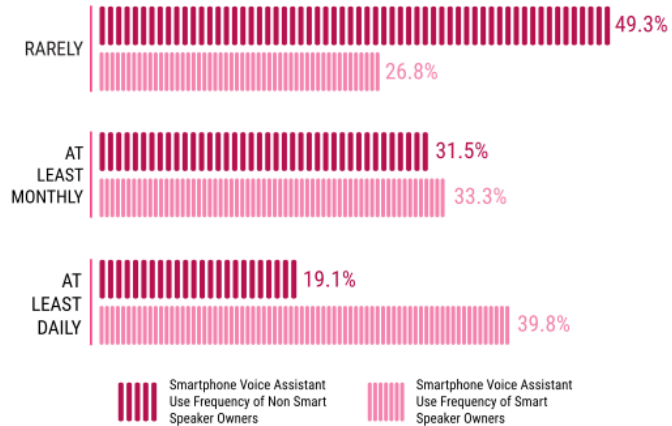
Smart Speaker Use Case Frequency January 2019



Per Household - U.S.



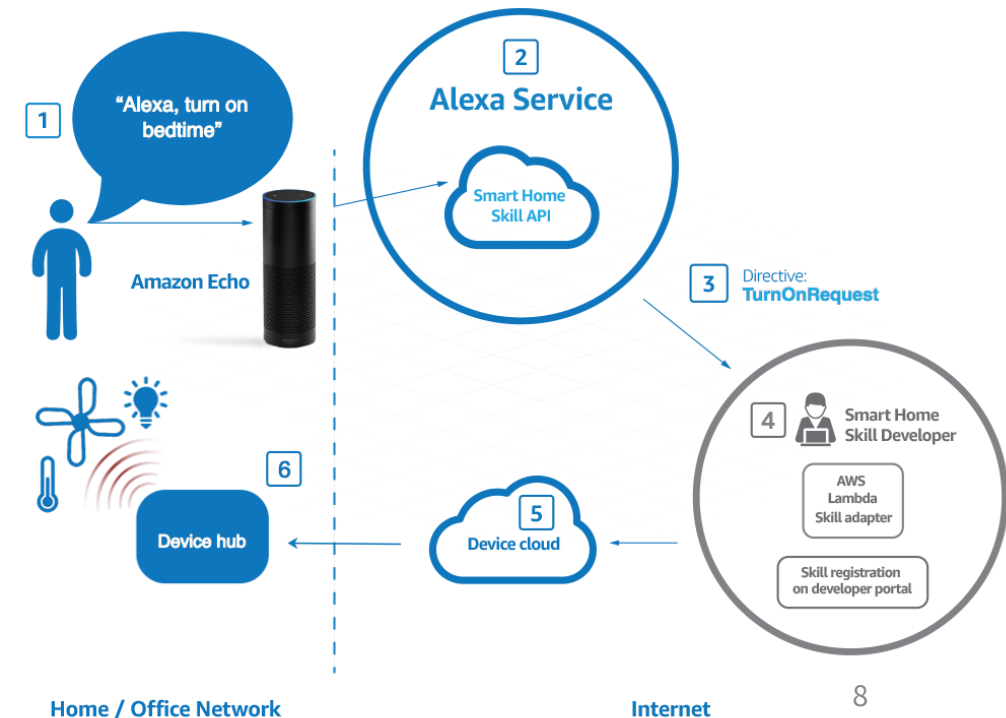
Voice Assistant Use Frequency on Smartphones by Smart Speaker Ownership



Source: Voicebot Smart Speaker Consumer Adoption Report Jan 2019

# How they work

- Device listens for wake word
  - after the wake word, everything is processed in vendor's cloud service
  - raw audio is sent to vendor
  - follow-up mode – no wake word needed for follow-up questions (device listens for 5-10sec after replying)
  - privacy concerns
- Intents – designed for each domain
  - NLU trained on examples
  - DM + NLG handcrafted
  - extensible by 3<sup>rd</sup> parties (Skills/Apps)
- No incremental processing

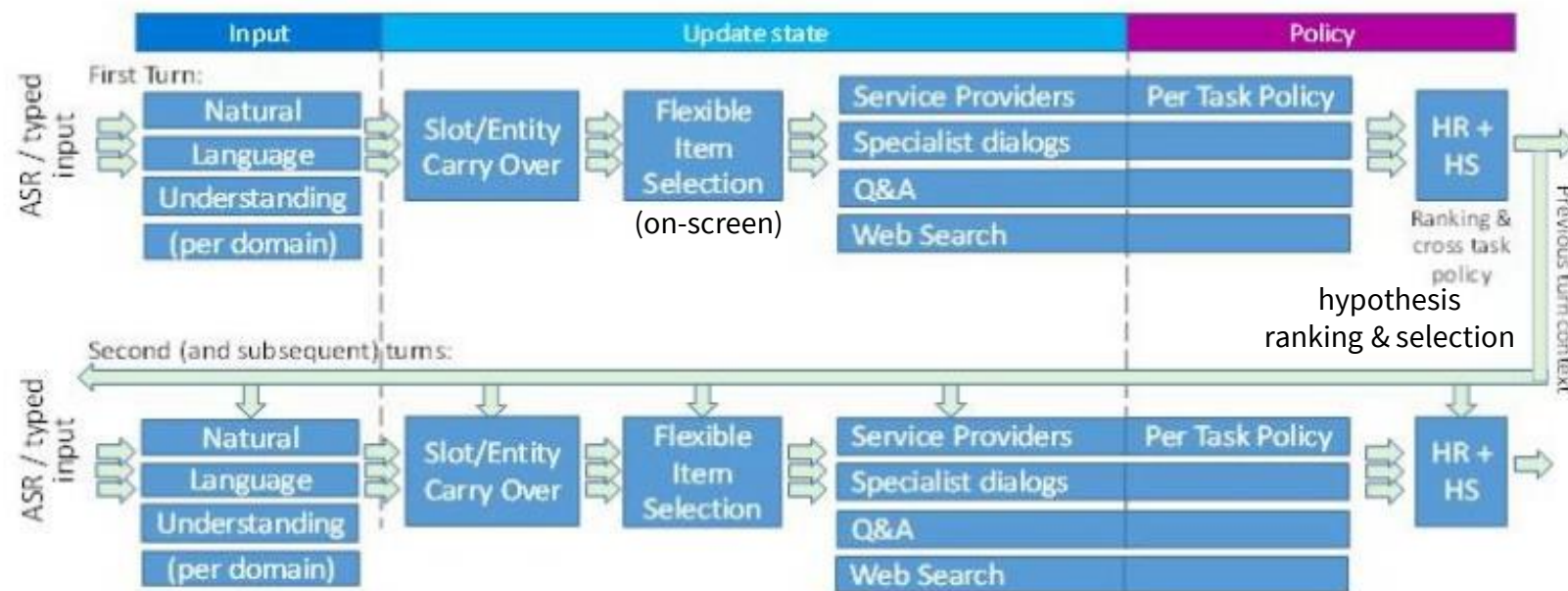




# How they work

- NLU includes domain detection
  - “web” domain as fallback
- Multiple NLU analyses (ambiguous domain), resolved in context
- State tracker & coreference
  - Rules on top of machine learning
  - All per-domain

Cortana structure



# Why they are cool

- ASR actually impressive
  - NLU often compensates for problems
- Range of tasks is wide & useful
- 1<sup>st</sup> really large-scale dialogue system deployment ever
  - not just a novelty
  - actually boosts voice usage in other areas (phone, car etc.)

Assistants & Accents

<https://youtu.be/gNx0huL9qsQ?t=41>



# Why they are not so cool

- Still handcrafted to a large part
  - *conversational architects* are a thing now
- Not very dialogue-y
  - mostly just one turn, rarely more than a few
- Language limitations
  - only available in a few major languages (En, Zh, Jp, De, Es, Fr, Kr [...])
- ASR still struggling sometimes
  - noise + accents + kids
  - not that far-field
  - helped *a lot* by NLU / domain knowledge

<https://youtu.be/CYvFxs32zvQ?t=65>



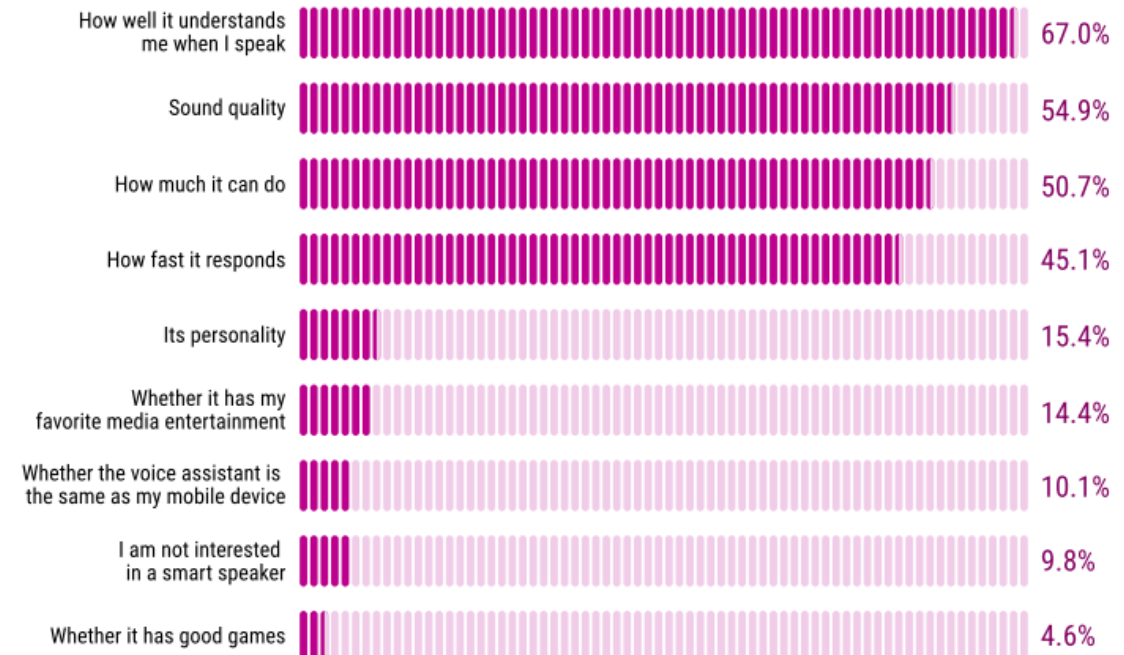
# Adding Skills/Apps

- Additional functionality by 3<sup>rd</sup> party developers
  - API/IDEs provided by vendors – see next time!
  - enabled on demand (similar to installing phone apps)
- Not 1<sup>st</sup>-class citizens
  - need to be invoked specially
    - *Alexa, tell Pizza Hut to place an order*
    - *Alexa, ask Uber to get me a car*
- There's thousands of them
  - many companies have a skill
  - many specific inventions
  - finance, fitness, food, games & trivia ...
  - much less used than the default ones

# What people care about in smart speakers

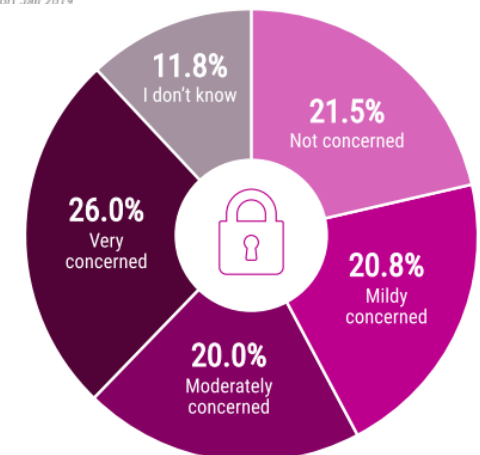
- **Understanding, features, speed**
  - personality / dialogue not so much
  - 3<sup>rd</sup> party apps not so popular (should work out-of-the-box)
  - commerce not so popular, but growing
- QA: music, news, movies
- Privacy concerns don't stop people from buying/using smart speakers
  - privacy-conscious 16% less likely to own one

What Qualities U.S. Users Value in Smart Speakers



Source: Voicebot Smart Speaker Consumer Adoption Report Jan 2019

U.S. Consumer Perception of Smart Speaker Privacy Risk



Source: Voicebot Smart Speaker Consumer Adoption Report Jan 2019



# Question answering

- integral & important part of assistants
  - broadest domain available, apart from web search
- QA vs. web search
  - QA needs a specific, unambiguous answer, typically a (named) entity
    - person, object, location [...]
    - ~ **factoid questions**
  - Needs to be within inference capabilities of the system

*Who is the president of Germany?  
How high is the Empire State Building?*

x

*Who is the best rapper?  
Who will become the next U.S. president?  
How much faster is a cheetah than an elephant?*

# Web search

- Given query, find best-matching **documents**
  - Over unstructured/semi-structured data (e.g. HTML)
- Basic search
  - Candidates: find matching word occurrences in index
  - Reranking: many features
    - Location of words (body, title, links)
    - Frequency of words (TF-IDF)
    - Word proximity
    - PageRank – weighing links to documents/webpages (how many, from where)
  - 2<sup>nd</sup> level: personalized reranking
- Query reformulation & suggestion

# QA approaches

- **Information Retrieval**

- Basically improved web search
- IR + phrase extraction
  - getting not just relevant documents, but specific phrases within them

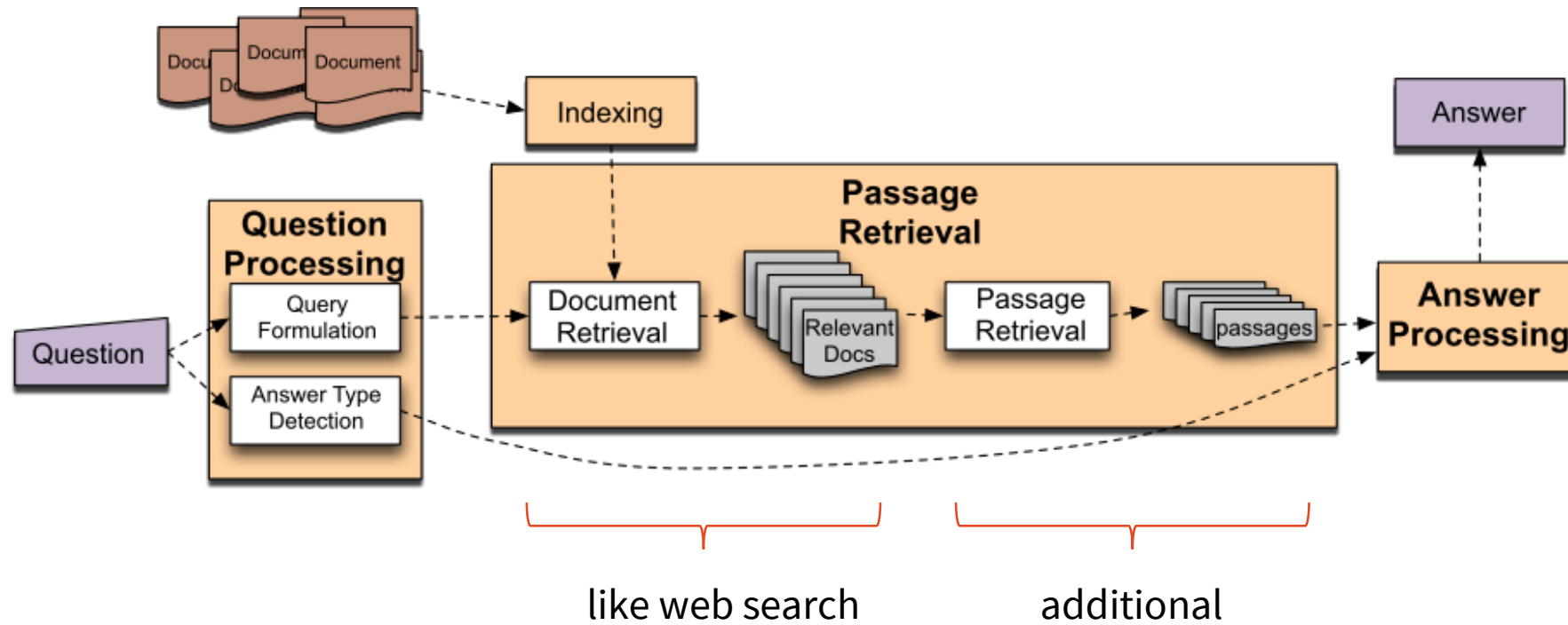
- **Knowledge Graphs**

- KGs – storage of *structured* information
  - 1) Semantic parsing of the query
  - 2) Mapping to KG(s)

- **Hybrid** (IBM Watson)

- candidates from IR
- reranking using KGs/semantic information

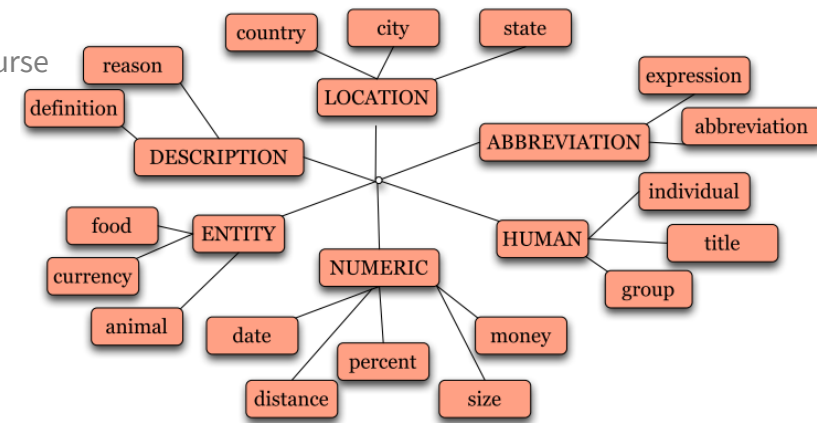
# IR-based QA Pipeline



from Jurafsky & Manning  
QA slides, Coursera NLP course

# Question Processing

Jurafsky & Manning  
QA slides, Coursera NLP course



- **Answer type detection**

- what kind of entity are we looking for?
- rules / machine learning (with rules as features)
- rules: regexes
  - headword = word right after wh-word

- **Named entity recognition**

- **IR Query formulation** – keyword selection

- ignore stop words (*the, a, in*)
- prioritize important words (named entities)
- stemming (remove inflection)

- **Question type classification** – definition, math...

- **Focus detection** – question words to replace with answer

- **Relation extraction** – relations between entities in question

- more for KGs, but can be used for ranking here

*Who is the [...] composer/football player [...]*  
*Which city is the largest [...]*



# IR Document Retrieval

- Candidates – find matching words in index (same as web search)
- Weighting
  - Frequency: **TF-IDF (term frequency-inverse document frequency)**
    - TF – document more relevant if term is frequent in it
    - IDF – document more relevant if term only appears in few other documents

$$\text{TFIDF} = (1 + \log f_{t,d}) \cdot \log \frac{N}{n_t}$$

← TF (log-scaled)
← # times  $t$  appears in  $d$ 
← IDF

← total # of documents  
 ← # of documents containing  $t$

- this is just one of many variants
- Other metrics – **BM25** – more advanced smoothing, heeds document length
- Proximity: also using n-grams in place of words

# IR Passage Retrieval

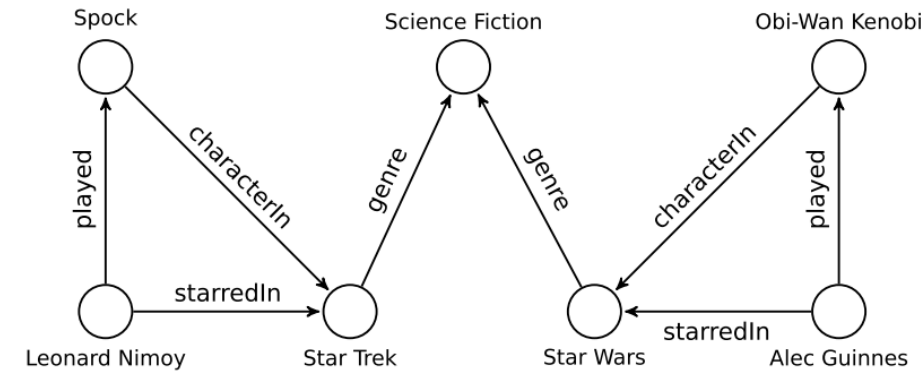
- Passage **segmentation** – split document into ~paragraphs
  - anything short enough will do
- Passage **ranking** – typically machine learning based on:
  - named entities & their type (matching answer type?)
  - # query words contained
  - query words proximity
  - rank of the document containing passage

# IR Answer Extraction

- **NER on passages** – looking for the right answer type
- 1 entity found → done
- More entities present → needs **another ranking**, based on:
  - answer type match
  - distance from query keywords in passage
  - novelty factor – not contained in query
  - position in sentence
  - semantic parse / relation
  - passage source rank/reliability

# Knowledge Graphs

- Large repositories of **structured, linked** information
  - **entities** (nodes) + **relations** (edges)
  - typed (for both)
  - entity/relation types form an **ontology** (itself a similar graph)
- Open KGs (millions of entities, billions of relations)
  - Freebase (freely editable, many sources, bought by Google & shut down)
  - DBpedia (based on Wikipedia)
  - Wikidata (part of Wikipedia project, freely editable)
  - Yago (Wikipedia + WordNet + GeoNames)
  - NELL (learning from raw texts)
- Commercial KGs: Google KG, Microsoft Satori, Facebook Entity Graph
  - domain specific: Amazon products, Domino's pizza [...]



from Jens Lehman's QA keynote

# RDF Representation

- RDF = Resource Description Framework
  - Most popular KG representation
  - Wikidata – different format but accessible as RDF
- **Triples:** <subject, predicate, object>
  - predicate = relation
  - subject, object = entities
  - can also include relation confidence (if extracted automatically)
- Entities & relations typically represented by URI (not always)
  - objects can also be constants (string, number)

subject: *Leonard Nimoy*  
 predicate: *played*  
 object: *Spock*  
 [confidence: 0.993]



# SPARQL

- Query language over RDF databases
  - relatively efficient
  - can query multiple connected triples (via ?variables)
- can be used directly
  - if you know the domain/application
- QA – need to map user question to this
  - or use IR-based methods instead

Wikidata: largest cities with female mayors

<https://query.wikidata.org/>

```

SELECT DISTINCT ?city ?cityLabel ?mayor ?mayorLabel
WHERE
{
  BIND(wd:Q6581072 AS ?sex)
  BIND(wd:Q515 AS ?c)

  ?city wdt:P31/wdt:P279* ?c . # find instances of subclasses of city
  ?city p:P6 ?statement . # with a P6 (head of government) statement
  ?statement ps:P6 ?mayor . # ... that has the value ?mayor
  ?mayor wdt:P21 ?sex . # ... where the ?mayor has P21 (sex or gender) female
  FILTER NOT EXISTS { ?statement pq:P582 ?x } # ... but the statement has no P582 (end date) qualifier

  # Now select the population value of the ?city
  # (wdt: properties use only statements of "preferred" rank if any, usually meaning "current population")
  ?city wdt:P1082 ?population .
  # Optionally, find English labels for city and mayor:
  SERVICE wikibase:label {
    bd:serviceParam wikibase:language "[AUTO_LANGUAGE],en" .
  }
}
ORDER BY DESC(?population)
LIMIT 10

```

# KG Retrieval

*How fast do jaguars run?  
What is a top speed of a jaguar?*

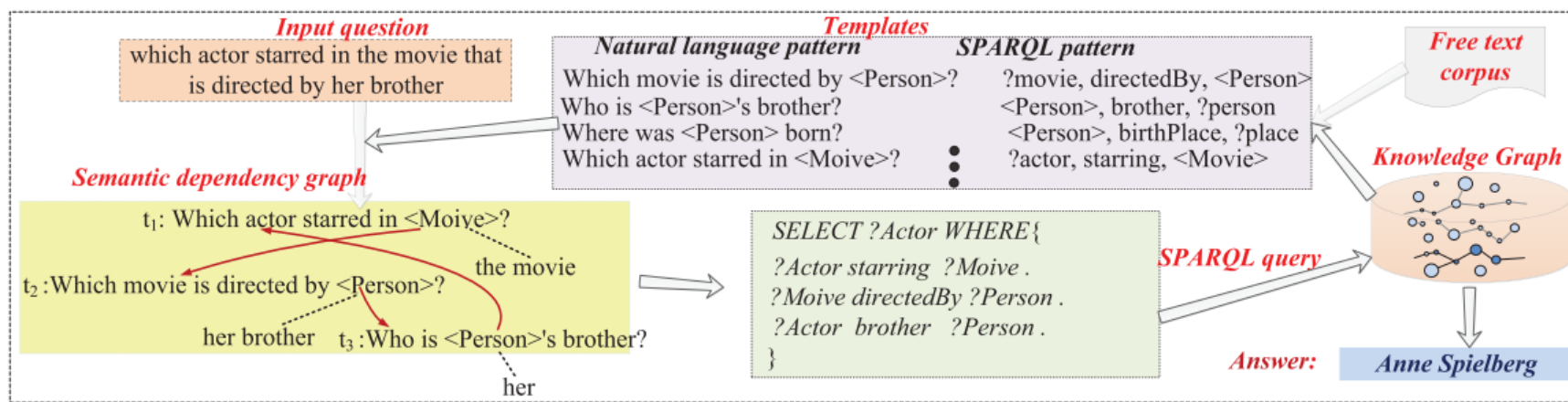
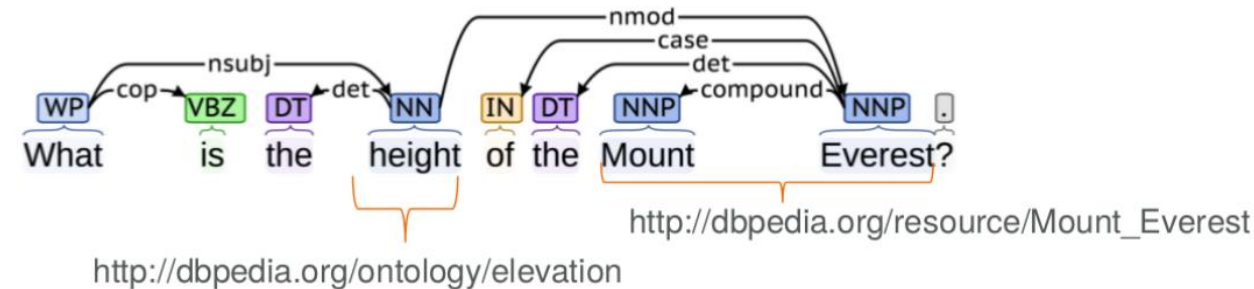
- Problem: **synonymy** – many ways to ask the same question
  - RDF relations have a specific surface form (not just *wd:1234*)
  - needs normalization/lexical mapping/usage of synonyms
    - WordNet expansion
    - stemming/lemmatization
    - multiple labels for entities/reasons
    - string similarity/word embeddings
  
- Problem: **ambiguity**

*How fast is a Jaguar [I-Pace]?*

  - needs entity/relation disambiguation/grounding/linking (to KG-compatible URIs)
  - context used to disambiguate (neighbour words, syntax, parts-of-speech)
  - KG itself used – closest/semantically related entities

# KG Retrieval

- **Semantic parsing** can be used for query normalization
- Dependencies help decompose complex questions
  - Doesn't have to be syntactic dependencies
  - Template mapping: map simple question patterns that have SPARQL equivalents



# KG Maintenance

- Information needs to be up-to-date
- Deduplication
- Ontology changes
  - need to version ontologies (and data)  
(for new/split/merged entity & relation types)
- Integrating multiple KGs
  - larger world knowledge coverage
  - company suppliers, mergers
  - → ontology bridging/mapping needed



*"Basically, we're all trying to say the same thing."*

<http://dit.unitn.it/~accord/RelatedWork/Matching/Noy-MappingAlignment-SSSW-05.pdf>

from Alex Marin's KG QA slides

# Ontology mapping

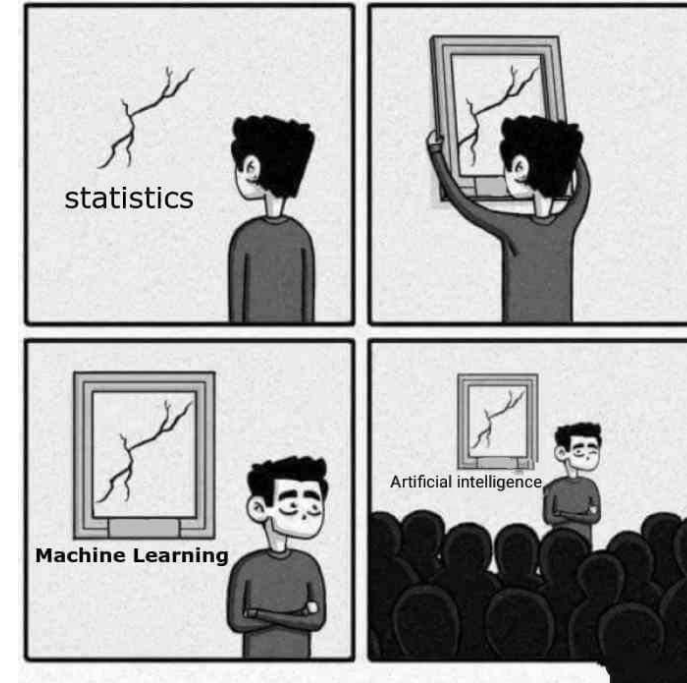
- Mismatch types
  - different labels (easiest)
  - same term, different thing & vice-versa
  - different modelling approaches (e.g. subclass or property?)
  - different granularity (more/less subclasses)
- Mappings
  - handcrafted (best results, but expensive)
  - rule-based – map into a common ontology
    - string distances, WordNet
  - graph-based – compare ontology structure
  - machine learning



# Machine Learning (Grossly Oversimplified)

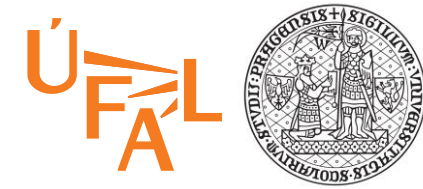
ML is basically function approximation

- function: data (**features**) → **labels**
  - discrete labels = **classification**
  - continuous labels = **regression**
- function shape
  - this is where different algorithms differ
  - neural nets: complex functions, composed of simple building blocks (linear, sigmoid, tanh...)
- **training/learning** = adjusting function parameters to minimize error
  - **supervised** learning = based on data + labels given in advance
  - **reinforcement** learning = based on exploration & rewards given online

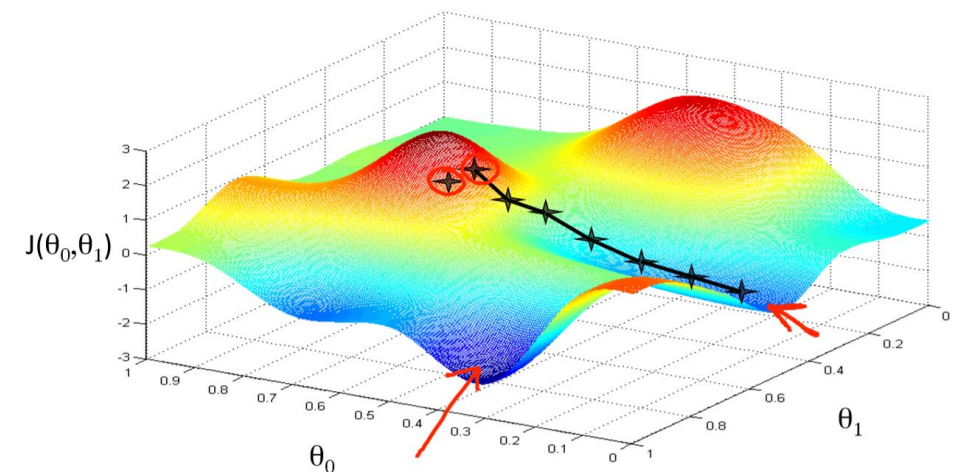


<https://towardsdatascience.com/no-machine-learning-is-not-just-glorified-statistics-26d3952234e3>

# Machine Learning (Grossly Oversimplified)



- training– **gradient descent** methods
  - minimizing a **cost function**  
(notion of error – given system output, how far off are we?)
  - calculus: derivative = steepness/slope
  - follow the slope to find the minimum – derivative gives the direction
  - **learning rate** = how fast do we go (needs to be tuned)
- gradient typically computed over **mini-batches**
  - random bunches of a few training instances
  - not as erratic as using just 1 instance, not so slow as computing over whole data
  - **stochastic gradient descent**
  - improvements: AdaGrad, Adam [...]
    - cleverly adjusting the learning rate



# Summary

- Virtual assistants/smart speakers are booming
  - large variety of tasks, interconnected
  - impressive ASR
  - typically handcrafted dialogue policy, NLG
- Question answering – factoids
  - a large part of assistants’ appeal, useful if integrated with tasks
  - IR approaches: word-based document retrieval, passage extraction, ranking
  - KG approach: semantic parsing & mapping to SPARQL queries
- Machine learning
  - finding the right function parameters by following cost function gradients
  - have a look at <http://jalammar.github.io/visual-interactive-guide-basics-neural-networks/>

# Thanks

## Contact me:

[odusek@ufal.mff.cuni.cz](mailto:odusek@ufal.mff.cuni.cz)  
room 424 (but email me first)

**Labs tomorrow**  
**9:00 SU1**

## Get the slides here:

<http://ufal.cz/npfl123>

## References/Further:

- Dan Jurafsky & Chris Manning's slides at Stanford/Coursera: <https://web.stanford.edu/~jurafsky/NLPCourseraSlides.html>
- Alex Marin's slides at Uni Washington: [https://hao-fang.github.io/ee596\\_spr2018/](https://hao-fang.github.io/ee596_spr2018/)
- Anton Leuski's slides at UCSC: <http://projects.ict.usc.edu/nld/cs599s13/>
- VoiceBot smart speaker report: <https://voicebot.ai/smart-speaker-consumer-adoption-report-2019/>
- Wikipedia pages of the individual KGs, assistants + [Smart speaker](#), [Okapi BM25](#), [TF-IDF](#)