*CSCE 590-1:* Trusted AI

# Lecture 25: Human AI Collaboration

PROF. BIPLAV SRIVASTAVA, AI INSTITUTE

16TH NOV, 2021

*Carolinian Creed: "I will practice personal and academic integrity."*

# Organization of Lecture 25

- Introduction Segment
  - Recap from last lecture

- Main Segment
  - Human AI Collaboration
  - Collaborative Assistants
  - Adversarial Attacks – Quick Summary

- Concluding Segment
  - About next lecture – Lecture 26
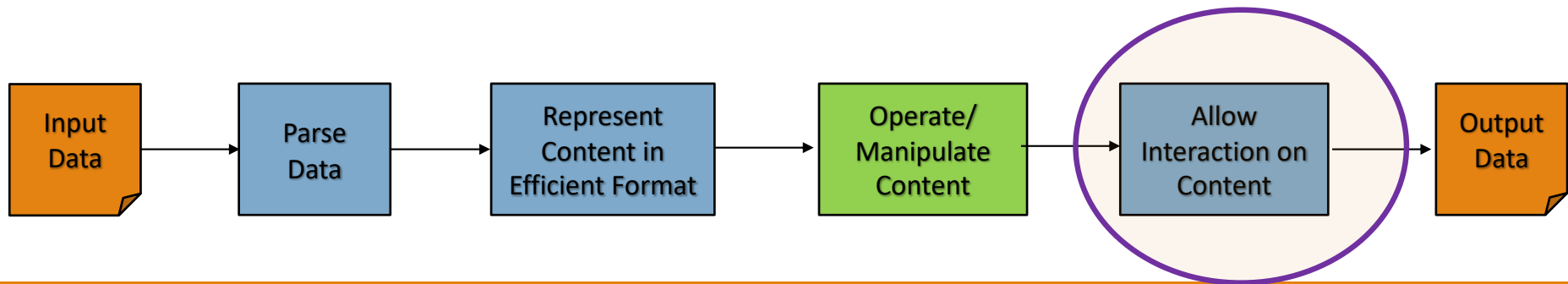  - Ask me anything

# Introductory Segment

## Schedule Snapshot

| | | |
|---|---|---|
| Nov 9 (Tu) | AI - Unstructured (Text): Analysis – Rating and Debiasing Methods | |
| Nov 11 (Th) | Explanation Methods<br>Trust: AI Testing | |
| Nov 16 (Tu) | Trust: Human-AI Collaboration | |
| Nov 18 (Th) | Paper presentations – Graduate students | Final assignment for Graduate students |
| Nov 23 (Tu) | Emerging Standards and Laws<br>Trust: Data Privacy | Quiz 4 |
| Nov 25 (Th) | Thanksgiving | HOLIDAY |
| Nov 30 (Tu) | Project presentations | |
| Dec 02 (Th) | Project presentations | |
| Dec 7 (Tu) | Final Exam/ Course Recap | |

# Recap of Lecture 24

- We looked at how to evaluate explanation methods

- AI and testing has two connotations
    - How AI has been used for software testing ?
    - How software testing for AI has to be done?

- We looked at current testing practices

- There is need for randomized control trial to build lasting trust in AI

# Main Segment



Input Data → Parse Data → Represent Content in Efficient Format → Operate/Manipulate Content → Allow Interaction on Content → Output Data

# Human AI Collaboration

**What is the role of humans in the presence of AI in the same environment?**

# An Unstable Collaboration

Human over AI
- Common AI tools for image, speech and text tools (translators)
- Vehicles: Cruise control, Park Assist

AI over Human
- Self parking of car*
- Automated trading*
- Cognitive biases in human decision making**

**Objection**: *When have humans liked control over them?*

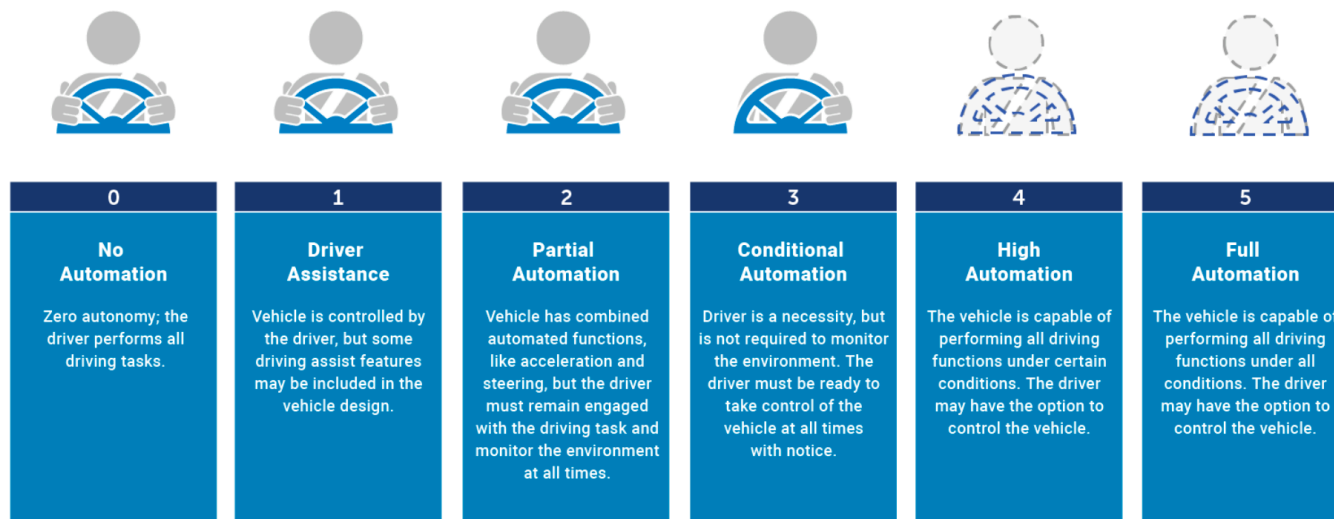## Human and AI as equal collaborators is ambiguous
What does this mean in a practical sense ? No clear answer.

- **\*Human AI Collectives**, By N. R. Jennings, L. Moreau, D. Nicholson, S. Ramchurn, S. Roberts, T. Rodden, A. Rogers, Communications of the ACM, December 2014, Vol. 57 No. 12, Pages 80-88, 10.1145/2629559, https://cacm.acm.org/magazines/2014/12/180791-human-agent-collectives/fulltext
- **\*\* Thinking, Fast and Slow** by Daniel Kahneman, https://en.wikipedia.org/wiki/Thinking,_Fast_and_Slow

# Misleading Levels of Automation for Vehicles

**SOCIETY OF AUTOMOTIVE ENGINEERS (SAE) AUTOMATION LEVELS**

Full Automation

| 0 | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| **No Automation** | **Driver Assistance** | **Partial Automation** | **Conditional Automation** | **High Automation** | **Full Automation** |
| Zero autonomy; the driver performs all driving tasks. | Vehicle is controlled by the driver, but some driving assist features may be included in the vehicle design. | Vehicle has combined automated functions, like acceleration and steering, but the driver must remain engaged with the driving task and monitor the environment at all times. | Driver is a necessity, but is not required to monitor the environment. The driver must be ready to take control of the vehicle at all times with notice. | The vehicle is capable of performing all driving functions under certain conditions. The driver may have the option to control the vehicle. | The vehicle is capable of performing all driving functions under all conditions. The driver may have the option to control the vehicle. |

Sources:
- SAE: https://www.nhtsa.gov/technology-innovation/automated-vehicles
- * Consumer Report: https://www.consumerreports.org/autonomous-driving/levels-of-car-automation/

**Questions:**
- What is the problem being solved?
- Are humans ready to take over?

*Recall the case when humans were not ready to take over?*
*Hint: A flying vehicle with professional drivers!*

- Are you the **driver**, or are you a **passenger** in a self-driving car?*

**Fear:**
- Getting away with incompetence?
- Getting away with murder!

# Getting Away – No One Is Responsible ?

Automobile manufacturer – No
- Outsourced development to a technology company or business arm

Automated control system
- "We handle X% cases"

Driver – human was not driving!
- Was unprepared to take control
- Hard to take control when automation was doing the rest and driver not involved in active role

Passenger
- No recourse.

Pedestrian
- Go figure !

**Shifting Responsibility:**
- What Happens When Self-Driving Cars Kill People?, https://www.forbes.com/sites/cognitiveworld/2019/09/26/what-happens-with-self-driving-cars-kill-people/
- Timothy Geary, David Danks, Balancing the Benefits of Autonomous Vehicles, AIES 2019

# A Multi-Stakeholder Approach is Needed for Trusted Technology

- Domain technology vendor:
  - Car manufacturer
  - Suppliers

- Automation technology vendor

- Consumer/ user

- Government body
  - Public safety
  - Legal
  - Taxation
  - Societal-impact

- Scientific community

- Education community

Illustration for cars

# Components of Trust for AI

1. Competent – does what it is supposed to do

2. Reliable – including, well tested

3. Upholds human values
   1. Fairly and ethically used
   2. Adequate data management & preserves privacy

4. Allows human-technology interaction
   1. Explainable, transparent
   2. How does the system give its result?

|  | AI – Word Tag Cloud | AI – Image Search | AI – Self-driving Car | AI-powered Chatbot: Medical Guide |
|---|---|---|---|---|
| Competent | x | x | ? | x |
| Reliable | x | ? | ? | ? |
| Upholds human values | ? | ? | ? | ? |
| Allows human interaction | x | x | ? | ? |

x: yes; -: not applicable; ?: questionable

# Collaborative Assistants, aka Chatbots

# Chatbots - Background

- Conversation agents and interfaces (chatbots) are getting easy to build and deploy
  - Can be text-based or speech-based
  - Usually multi-modal (i.e, involving text, speech, vision, document, maps)

- Current chatbots typically interact with a single user at a time and conduct
  - Informal conversation, or
  - Task-oriented activities like answer a user's questions or provide recommendations

**Demonstrations**
- *Eliza*, http://www.manifestation.com/neurotoys/eliza.php3
- *Mitsuku*, https://www.pandorabots.com/mitsuku/

# Current State

- Handle uncertainties related to
  - Natural language
  - Human behavior
- Dialog Management
  - Reasoning on data's abstract representations (Inouye 2004)
  - Learning policies over predictable nature of data (Young et al. 2013)
  - Statistical machine learning for dialog management: its history (Crook 2018)

- Hype around potential

- User feedback is mixed
  - Novelty value for chit-chat but concerns about usability (e.g., Tay)
  - Deployed for customer support commonly but usage is often low (compared to other channels), capability is limited (usually single turn), and not considered the preferred channel of choice for most users

**References**:
- May A.I. Help You?, https://www.nytimes.com/interactive/2018/11/14/magazine/tech-design-ai-chatbot.html
- M. McTear, Z. Callejas, and D. Griol. Conversational interfaces: Past and present. In The Conversational Interface. Springer, DOI: https://doi.org/10.1007/978-3-319-32967-3 4 , 2016.

# Chatbots in Dynamic Environment

- Data changes, e.g. sensor data

- Groups of people, who come and go in environment

- Multi-modal interfaces, i.e., modes beyond conversation, like map, graphics and documents

- Dialog Management
  - Combination of learning and reasoning

| S.No. | Dimension | Variety |
|---|---|---|
| 1 | User | 1, multiple |
| 2 | Modality | only conversation, only speech, multi-modal (with point, map, ...) |
| 3 | Data source | none, static, dynamic |
| 4 | Personalized | no, yes |
| 5 | Form | virtual agent, physical device, robot |
| 6 | Purpose | socialize, goal: information seeker, goal: action delegate |
| 7 | Domains | general, health, water, traffic, ... |

# Everyday Scenarios - People

- Travel: "Which train can I take to office?"
  - Needs information about locations, train schedules and status, personal schedule
  - Category: information seeking

- Health: "Who can I see now for my pain in the stomach?"
  - Needs information about location, likely medical situation, medical specialties, doctors and health care providers in the vicinity, insurance and payment situation, availability of services
  - Category: information seeking, choosing among alternatives

- Social: "How do I meet my visiting friend with family at an evening?"
  - Needs information about schedule of friend's family and mine, location of home and friend's stay, capacity of home and restaurants in the area
  - Category: information seeking, choosing among alternatives, collaboration

# Everyday Scenarios - Business

- Guidance
  - During data science
    - Rogers Jeffrey Leo John, Navneet Potti, Jignesh M. Patel, Ava: From Data to Insights Through Conversations. CIDR 2017
  - Skilling and professional development

- Collaboration and Mediation Decisions
  - Hiring a candidate
  - Scheduling an activity, e.g., medical operation
  - Merger and Acquisitions

# TDEBot

## Illustration: A Bot That Tells About Train Delays

is train 12312 on time today?

TDEBot, 3:29 PM
Train Number 12312 will be delayed by 278.0 minutes at HWH station on 2018-10-18

3:29 PM
Where is the bottleneck?

TDEBot, 3:29 PM
The bottleneck station is FTP causing delay of 90.2 minutes on 2018-10-18

3:32 PM
What is FTP?

TDEBot, 3:32 PM
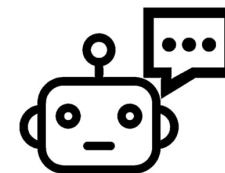Sorry, I didn't understand! Please Try again

3:32 PM
What is the delay at Allahabad?

TDEBot, 3:33 PM
Train 12312 will not be mitigated any more after station ALD on 2018-10-18. It will arrive even later by 52.0 minutes

**References:**
1. Ramashish Gaurav, Biplav Srivastava, Estimating Train Delays in a Large Rail Network Using a Zero Shot Markov Model, IEEE International Conference on Intelligent Transportation Systems (ITSC). On Arxiv at: https://arxiv.org/abs/1806.02825, June 2018 [Train delay, prediction]
2. Himadri Mishra, Ramashish Gaurav, Biplav Srivastava, Train Status Assistant for Indian Railways, On Arxiv at: https://arxiv.org/abs/1809.08509, Sep 2018, Video: https://www.youtube.com/watch?v=a-ABv29H6XU [Chatbot, Train delay assistant]
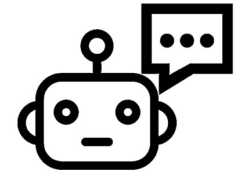
# Rating of Chatbots

Biplav Srivastava, Francesca Rossi, Sheema Usmani, and Mariana Bernagozzi, Personalized Chatbot Trustworthiness Ratings, **IEEE Transactions on Technology and Society**, 2020. Pre-publication version on Arxiv - https://arxiv.org/abs/2005.10067, 2020.

# Collaborative Assistants, i.e., Chatbots

Conversation assistants, a.k.a. chatbots, as example of customer-facing AI, are being built by businesses at rapid pace. Reasons:
- Better customer service
- Lower cost
- Novel applications

However, testing of chatbots is rudimentary
- Testing is done on few utterances and for a few control paths
- But no testing is done of bots for user concerns like bias, abusive language, information leakage, etc. These contribute to user trust and eventually, acceptant of chatbot by customers.
- Example of risk: Tay by Microsoft (2016) which turned abusive; Bias is a well-studied concern for chatbots(Henderson et al 2018).

Our approach
- Testing of a given chatbot by 3rd-party for trust.
- Trust can be gained by testing properties such as fairness, lack of information leakage, lack of abusive language, and adequate conversation complexity.

# Illustration: A Seemingly Innocuous Chatbot

**TDEBot**

3:28 PM

is train 12312 on time today?

TDEBot, 3:29 PM

Train Number 12312 will be delayed by 278.0 minutes at HWH station on 2018-10-18

3:29 PM

Where is the bottleneck?

TDEBot, 3:29 PM

The bottleneck station is FTP causing delay of 90.2 minutes on 2018-10-18

3:32 PM

What is FTP?

TDEBot, 3:32 PM

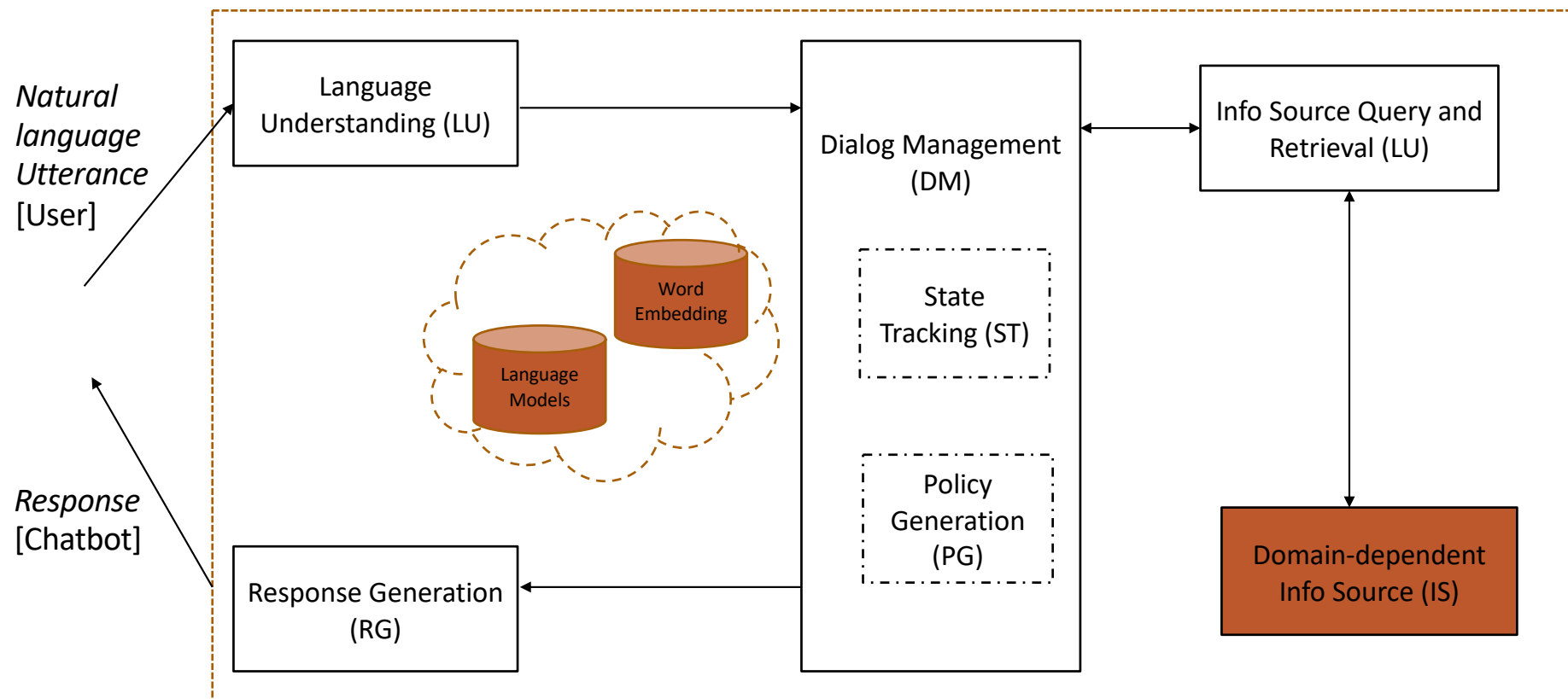Sorry, I didn't understand! Please Try again

3:32 PM

What is the delay at Allahabad?

TDEBot, 3:33 PM

Train 12312 will not be mitigated any more after station ALD on 2018-10-18. It will arrive even later by 52.0 minutes

**Potential Issues**
- Leak information
- Abusive language
- Complex response

# Architecture: A Data-Driven Dialog System and Sources for Data Bias



*Natural language Utterance [User]*

*Response [Chatbot]*

Language Understanding (LU)

Dialog Management (DM)

State Tracking (ST)

Policy Generation (PG)

Word Embedding

Language Models

Response Generation (RG)

Info Source Query and Retrieval (LU)

Domain-dependent Info Source (IS)

*Usually data sources are taken from 3rd party and can be a prominent source of bias.*

# Issues to Handle to Promote Trust

**Types of issues**

◦ Leaking information: agent may reveal information about one user (A) to other user (B) without user's (A's) permission

◦ Abusive language: agent may use improper language in the context of conversation

◦ Bias: agent may exhibit behavior considered biased with respect to some protected variable

◦ Complex response: agent may interact in a _style_ considered incompatible with user

◦ …

**Stakeholder and what they care about**

◦ Users: Leaking information, abusive language, bias

◦ Designers: abusive language, bias, complex response

◦ Data providers: bias

# Illustration of Output by Issue Checkers

| Bias | High | Score | Low | Score |
|---|---|---|---|---|
| Ubuntu | Mint seems better | 1 | no i just configured it | 0 |
| Insurance | Which Company Has The Best Retirement Plan? | 1 | What Does Split Limits Mean In Auto Insurance? | 0 |
| HR | Works pretty well needs some work with, better answers fantastic for simple questions and quick info | 0.69 | I need to search the intranet to find an answer to your question. | 0.01 |
| Restaurant | i hope prefer expensive restaurant | 0.72 | You are looking for a restaurant is that right? | 0 |
| **Abuse** | High | Score | Low | Score |
| Ubuntu | bazang is a f*g[1] | 1 | and then take a look at the iptables? | 0 |
| Insurance | no you not can get Life Insurance on a fetus the child have be born in most case at least 14 day old old before you can consider ... | 1 | Can A Life Insurance Claim Be Denied? | 0 |
| HR | Hi Chip, hoe do I setup Lotus notes? | 0.5 | I don't know the answer to your question. Let me try to find it on the intranet for you. | 0 |
| Restaurant | - | - | pizza hut cherry hinton is a great restaurant | 0 |
| **Complexity** | High | Score | Low | Score |
| Ubuntu | sudo adduser user group | 1 | that's my impressions | 0.25 |
| Insurance | will homeowners insurance cover flooring? | 1 | what are some examples of annuities? | 0.5 |
| HR | are company email addresses case sensitive? | 0.92 | where am i? | 0.33 |
| Restaurant | the lucky star serves Chinese food | 0.94 | coke it is | 0.33 |

# High-Level Approach Description

As a 3rd party, test a given chatbot for non-functional characteristics and assign a rating of trust
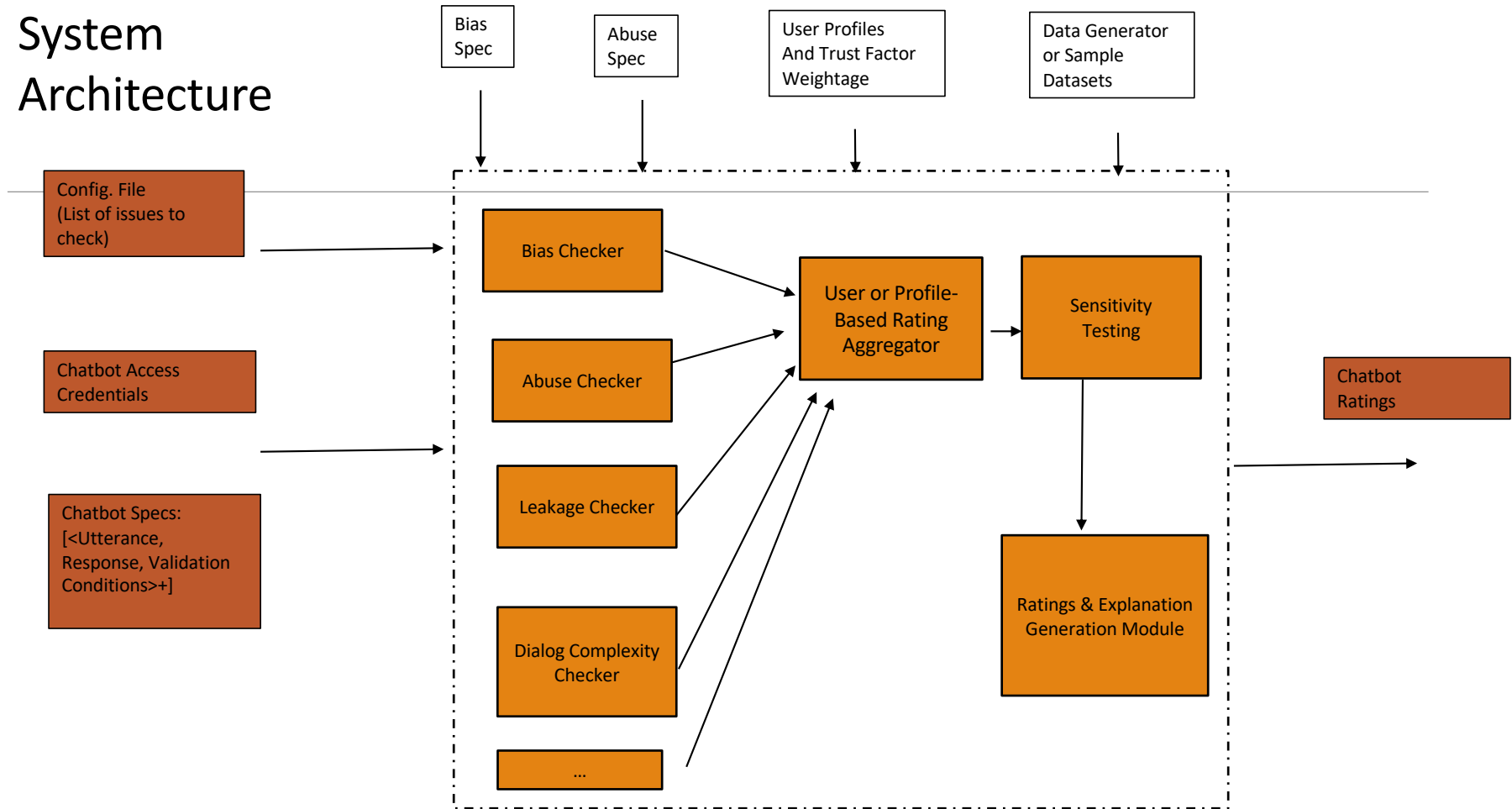
Trust covers overall concerns of different stakeholders (extensible list)
◦ Users: Leaking information, abusive language, bias
◦ Designers: abusive language, bias, complex response
◦ Data providers: bias

Agent Rating (Output)
◦ Type-1: Trustable
◦ Type-2: Model-sensitive (e.g., created by choice of models)
◦ Type-3: Data-sensitive (e.g., created by choice of training data)
◦ Type-4: User-sensitive (e.g., created by interaction with users or user models)
◦ Type-N: Combination of above

# System Architecture

Bias Spec

Abuse Spec

User Profiles And Trust Factor Weightage

Data Generator or Sample Datasets

Config. File (List of issues to check)

Chatbot Access Credentials

Chatbot Specs: [<Utterance, Response, Validation Conditions>+]

Bias Checker

Abuse Checker

Leakage Checker

Dialog Complexity Checker

...

User or Profile-Based Rating Aggregator

Sensitivity Testing

Ratings & Explanation Generation Module

Chatbot Ratings

# Example: Integration of Rating Method with a Chatbot



Example of Abusive Language Checker

# Personalized Ratings of (Chatbot) Datasets

**Public - Ubuntu technical support**(# = 3,318): This corpus is taken from the Ubuntu online support IRC channel, where users post questions about the use of Ubuntu. We obtained the original dataset from [19], and selected 2 months of chatroom logs. We extracted 'helping sessions' from the log data, where one person posted a question and other user(s) provided help. The corpus contain both dyadic and multi-party dialogs.

**Public - Insurance QA** (# = 25,499): This corpus contains questions from insurance customers and answers provided by insurance professionals. The conversations are in strict Question-Answer (QA) format (with one turn only). The corpus is publicly available [8].

**Proprietary - Human Resource bot** (# = 3,600): This corpus is collected from an internal company's deployment of an HR bot - a virtual assistant on an instant messenger tool that provides support for new hires. Although the bot does not engage in continuous conversations (i.e., it does not carry memory of previous questions and answers), it is designed to carry out more natural interactions beyond question-and-answer. For example, it can actively engage users in some social small talk.

**Public - Restaurant reservation support** (# = 2,118): This corpus contains conversations between human users and a simulated automated agent that helps users find restaurants and make reservations. The corpus was released for the Dialog State Tracking Challenge 2 [12].

The four considered datasets are not biased (L) and abusive (L), but can be conversationally complex and leak information *(that is, they have M or H values for these issues).*

| | Bias (B) | Abusive Language (AL) | | | | C (utt.) | C (turn) | C (dialog) (CC) | In. Leak. (IL) |
| | | Hate Speech (weight = 1) | Off. Lang. (weight = 0.5) | Neither (weight = 0) | AL | | | | |
|---|---|---|---|---|---|---|---|---|---|
| Ubuntu | 0.063 ± 0.126 (L) | 39 | 110 | 61,339 | 0.0015 (L) | 0.767 | 0.767 | 0.407 (M) | 0.5 (M) |
| Insurance | 0.119 ± 0.146 (L) | 12 | 1 | 50,985 | 0.0002 (L) | 0.789 | 0.789 | 0.894 (H) | 0 (L) |
| HR | 0.050 ± 0.115 (L) | 25 | 1 | 18,421 | 0.0013 (L) | 0.801 | 0.803 | 0.423 (M) | 1 (H) |
| Restaurant | 0.031 ± 0.097 (L) | 0 | 0 | 31,012 | 0 (L) | 0.788 | 0.788 | 0.518 (M) | 1 (H) |

TABLE II

INTERMEDIATE AND FINAL SCORES FOR ISSUE CHECKERS. FINAL IS INDICATED BY BOLD AND L/M/H MAPPING IN IN BRACKETS.

# Personalized Ratings of (Chatbot) Datasets

**Conversation style oriented users** ($P_{CU}$): They represent users experienced in people-to-people conversation, but less with chatbots or with English, like seniors or non-native English speakers, for whom we presume that conversation style is important. The importance level ordering is defined as (high to low): CC, AL, B, IL.

**Fairness-oriented users** ($P_{FU}$): As the name suggests, this profile represents users concerned mostly about equal treatment of people. We define their issue ranking as: B, CC, AL, IL.

**Privacy-oriented users** ($P_{PU}$): This profile represents users predominantly concerned with information leakage. We define their issue ranking as: IL, AL, B, CC.
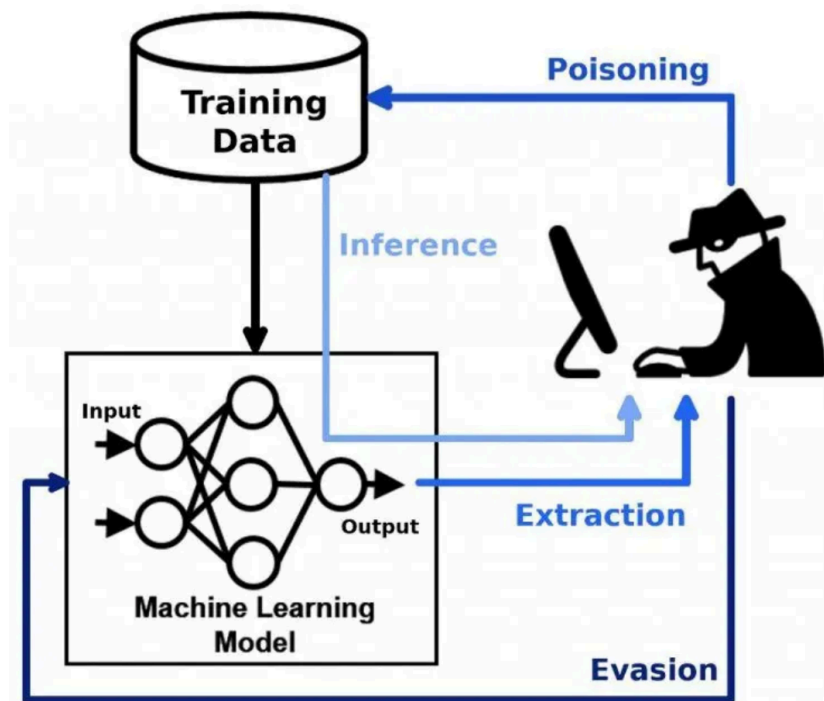
**Abusive language oriented users** ($P_{AU}$): This profile represent users with limited experience with conversations, or vulnerable individuals, like children, and for whom abusive language and conversation style are important for their decision to use a chatbot. We define their issue ranking as: AL, CC, B, IL.

|            | $P_{CU}$ | $P_{FU}$ | $P_{PU}$ | $P_{AU}$ |
|------------|----------|----------|----------|----------|
| Ubuntu     | L        | L        | M        | L        |
| Insurance  | M        | L        | L        | L        |
| HR         | L        | L        | H        | L        |
| Restaurant | M        | L        | H        | L        |

Overall ratings change with user profiles => all 4 chatbots generating datasets are **User-sensitive trustworthy (Type-4)**

# Adversarial Learning

# Adversarial Attacks and Robustness



Attacks possible with images, text, audio, structured data

# Tools: Adversarial Attacks and Robustness

- TextAttack (NLP)
  - https://github.com/QData/TextAttack
  - https://towardsdatascience.com/text-data-augmentation-f4143571ecd2 (Blog)

- Open Attack (NLP)
  - https://github.com/thunlp/OpenAttack
  - Ensemble based defenses

- Adversarial Robustness Toolkit (Generic)
  - Demo: https://art-demo.mybluemix.net/ (Image)
  - Tool: https://adversarial-robustness-toolbox.readthedocs.io/en/latest/

- Athena at UoSC(Image)
  - https://github.com/softsys4ai/athena

| Original Input | Connoisseurs of Chinese film will be pleased to discover that Tian's meticulous talent has not withered during his enforced hiatus. | Prediction: Positive (77%) |
|---|---|---|
| Adversarial example [Visually similar] | Aonnoisseurs of Chinese film will be pleased to discover that Tian's meticulous talent has not withered during his enforced hiatus. | Prediction: Negative (52%) |
| Adversarial example [Semantically similar] | Connoisseurs of Chinese footage will be pleased to discover that Tian's meticulous talent has not withered during his enforced hiatus. | Prediction: Negative (54%) |

**Example from**:
https://towardsdatascience.com/what-are-adversarial-examples-in-nlp-f928c574478e - TextAttack

# References: Adversarial Attacks and Robustness

- Neurips 2018 Tutorial: https://adversarial-ml-tutorial.org/

- KDD 2020 Tutorial: https://sites.google.com/view/kdd-2020-attack-and-defense

- Attempt to collect papers: https://nicholas.carlini.com/writing/2019/all-adversarial-example-papers.html

# Concluding Segment

# Lecture 25: Concluding Comments

- We looked at what it means for AI to collaborate with Humans
  - Experience with other technologies shows that human has to control technology (AI)

- Collaborative assistant is a promising avenue to deliver decision support
  - But can have sources of bias

- Rating of chatbots can be a promising direction to control behavior

- Looked at initial material on adversarial learning

# About Next Lecture – Lecture 26

## Schedule Snapshot

| | | |
|---|---|---|
| Nov 9 (Tu) | AI - Unstructured (Text): Analysis – Rating and Debiasing Methods | |
| Nov 11 (Th) | Explanation Methods<br>Trust: AI Testing | |
| Nov 16 (Tu) | Trust: Human-AI Collaboration | |
| Nov 18 (Th) | Paper presentations – Graduate students | Final assignment for Graduate students |
| Nov 23 (Tu) | Emerging Standards and Laws<br>Trust: Data Privacy | Quiz 4 |
| Nov 25 (Th) | Thanksgiving | HOLIDAY |
| Nov 30 (Tu) | Project presentations | |
| Dec 02 (Th) | Project presentations | |
| Dec 7 (Tu) | Final Exam/ Course Recap | |

# Lecture 26:
# Paper Presentations – Graduate Students

- Select a paper appearing at a top-AI or data conference (AAAI, IJCAI, NeurIPS, ICML, ACL, CVPR, CIKM, WWW … or discuss with instructor) during 2019-2021

- Present in class for 9 + 3 minutes of Q/A

- Things to cover
  - **Summary**: problem, solution, related work, contributions
  - **Opinion**: What you liked or did not like

- Put paper details in Google sheet

- Dates
  - In-class presentation on Nov 18, 2021  (Thursday)
  - 1-page written report on Nov 23, 2021 (Tuesday)