world**usability**day

# A Gold Mining Adventure

Using Natural Language Processing, Machine Learning, and Human-Centered Design to Find Gold in Unstructured Data

**Chris Schilstra, MBA, SAFe, CSPO, CSM**

Innovation Program Manager

Tantus Technologies, Inc.

November 12, 2020

Certified Section
508 Compliant

**Human-Centered Design**
Center of Excellence

# AGENDA

## Interpreting Unstructured Listening Session Data

- Project Goals and Objectives
- Listening Session Process
- Find Gold in Unstructured Data
- Optimize NLP Visualizations With HCD
- HCD Project Outcomes
- Additional HCD Considerations/Puzzles

**Human-Centered Design**
Center of Excellence

# PROJECT GOALS AND OBJECTIVES

Improve the review process of unstructured listening session data by leveraging Artificial Intelligence and Human-Centered Design.

**Goals:**

Use HCD, Machine Learning (ML) and Natural Language Processing (NLP) to:

➤ Achieve high reliability (>80%) accuracy for review of listening session data

➤ Minimize/manage opportunities for human bias

**Objectives:**

- Derive Themes
- Determine Sentiment

| Listening Session Topic | Data Set |
|---|---|
| **Prior Authorization** | Data from several transcribed, national listening sessions summarized into one data set |

**Human-Centered Design**
Center of Excellence

# LISTENING SESSION PROCESS

| Without AI | With AI |
|---|---|
| Conduct listening sessions with stakeholders | Conduct listening sessions with stakeholders |
| Transcribe recorded audio into text | Transcribe recorded audio into text |
| **Manually** sort text snippets into thematic groups | **Use NLP** to categorize data into thematic groups |
| **Manually** identify the sentiment of themes | **Use ML** to assign sentiment to themes |
| **Retype findings** into a management summary | Visualize themes/sentiment; **automate reporting** |

**Human-Centered Design**
Center of Excellence

CCSQ WORLD USABILITY DAY

# USE HCD AND NLP TO

Find Gold In Unstructured Data

**Human-Centered Design**
Center of Excellence

# FIND GOLD IN UNSTRUCTURED DATA

# THEMES

**Themes**

- Derive themes within unstructured data using Natural Language Processing
- Themes are scored and ranked by Tf-Idf* within the unstructured data

**Sub-themes**

- Derive sub-themes within themes through verb/adverb/adjective proximity and word frequency
- Sub-themes are scored and ranked by frequency

**\*Tf-Idf = Term frequency, Inverse document frequency** – an algorithm used to evaluate how relevant or important a term is within a data set. The importance of a word increases proportionally to the number of times that word appears in the data set but is offset by the frequency of the words in the data set.

See Wikipedia at https://en.wikipedia.org/wiki/Tf%E2%80%93idf.

# FIND GOLD IN UNSTRUCTURED DATA
# HCD FOR THEMES

- Replace manual derivation of themes, which is highly prone to bias, with NLP-driven, numerically-based derivation of themes

  ➤ Automates/accelerates the speed of extracting themes

  ➤ Reduces bias

**Human-Centered Design**
Center of Excellence

# FIND GOLD IN UNSTRUCTURED DATA
## SENTIMENT

Derive the Sentiment of unstructured data using machine learning instead of manual sentiment assignment, thereby reducing bias.

Positive

Negative

Neutral

**Human-Centered Design**
Center of Excellence

# FIND GOLD IN UNSTRUCTURED DATA
# USING ML TO FIND SENTIMENT

## Key

- **Build Process** (blue)
- **Execution Process** (green outline)
- **SME Support** (orange)

**Receive Data**

**Iterate model build/train until targeted accuracy achieved**

**Build Training Data**

**Train Model (75%)**

**Test Model (25%)**

**Approve and Deploy Model**

**SME**

SME Support:
Identify, validate, and correct domain-specific data.

**Receive and Transform Data**

**Execute ML Model**

**View/Monitor Results**

**SME**

SME Support:
Validate domain-specific execution results

**Human-Centered Design**
Center of Excellence

# FIND GOLD IN UNSTRUCTURED DATA
# AI MODEL ACCURACY

- Trained and executed the following models to compare results
  - BERT (Google)
  - RoBERTa (Facebook)
  - DistilBERT (A smaller and lighter version of the BERT model)
- A pre-trained model with 150 Gb of data
- Results – RoBERTa produced superior results

| Model | Creator | Accuracy | F1 Score | MCC | Eval_Loss |
|---|---|---|---|---|---|
| **RoBERTa** | **Facebook** | **0.8210** | **0.8210** | **0.6746** | **0.5046** |
| DistilBERT | Cornell University | 0.7354 | 0.7354 | 0.4932 | 0.7595 |
| BERT | Google | 0.7310 | 0.7310 | 0.6282 | 0.5615 |

**Human-Centered Design**
Center of Excellence

CCSQ WORLD USABILITY DAY

# FIND GOLD IN UNSTRUCTURED DATA
## LISTENING SESSION THEMES WITH SENTIMENT



**Total Comments = 2,241**

Sentiment
- Neutral
- Negative
- Positive

**Comment Sentiment Overview**

- 515 Positive
- 1,086 Negative
- 640 Neutral

**Themes / Sub Themes**

Themes:
- prior authorization
- prior auth
- prior authorization process
- medicare advantage
- health plan
- insurance company
- patient care
- administrative burden
- medical necessity
- medicare advantage plan
- authorization request
- pre authorization
- prior authorization request
- step therapy
- real time
- health care
- utilization management
- clinical documentation
- pre auth
- prior authorization form
- patient access
- home health
- d plan
- acute care
- emergency room
- electronic prior authorization
- durable medical equipment
- prior authorization list
- physical therapy
- prior authorization approval

Total Sentences (0, 100, 200, 300, 400, 500, 600, 700)

# PROJECT GOALS AND OBJECTIVES ATTAINED

**Goals achieved:**

- Achieved high reliability (>80%) accuracy for sentiment assignment of listening session data
- Minimized human bias by utilizing ML and NLP modeling

| Objectives | Results |
|---|---|
| **Derive Themes/sub-themes** | • Deep learning algorithms determine the themes based on prevalence of text content and determine sub-themes based on theme proximity and frequency |
| **Determine Sentiment** | • Achieved >80% accuracy<br>• Machine learning algorithms determine sentiment at multiple levels:<br>  • Overall sentiment<br>  • Sentiment by themes<br>  • Sentiment by stakeholders |

# Optimize NLP Visualizations With HCD

**Human-Centered Design**
Center of Excellence

# OPTIMIZING NLP VISUALIZATIONS – ATTEMPT 1

## Themes and Sub-themes Tree View (Sentence-driven)

**Sentence**

**Sub-theme**

**Theme**

establishes a prior authorization program , it should create a process or provide staffing for hospitals

our providers are not able

to

cms and ma plans can require those providers

receive

mha urges cms to use its authority to ensure hospitals are paid for services that

we believe cms and ma plans should only require outlier billers

they -- we get physicians to sign off on it , but they receive that

they -- the company won't release it unless it's -- they receive the

if an insurer establishes a

physician anesthesiologists often encounter unpaid patient medical bills in such cases where the patient has not received

prior **authorization**

, and they see it through .

program , it should create a process or provide staffing for hospitals to receive prior authorization 24 hours a day , 365 days a year .

24 hours a day , 365 days a year .

until the day of the scheduled procedure leaving no time to adequately counsel the patient on the cost of the service .

but underwent the surgical procedure .

while sparing compliant providers from this added burden .

# OPTIMIZING NLP VISUALIZATIONS – ATTEMPT 2

Themes and Sub-themes (Theme-driven View)



**Sentence**

**Sub-theme**

**Theme**

of breast cancer required prior authorization.

**require**
Context Sentences: 65

theres no advantage whatsoeverinformation concerning a medicare supplement beneficiarys hospitalization is not readily available to their plans, whether the beneficiary comes to the hospital through the emergency department or

**obtain**
Context Sentences: 38

however, for the most part, the challenges that were arising were not on supplies that required prior authorization, not for the large part.

**need**
Context Sentences: 18

will chemotherapy drugs require prior authorization in the near future?

**receive**
Context Sentences: 13

instead of requiring prior authorization for these procedures, tests and medications up front, if the carrier felt it necessary, they could monitor usage on the backend to look for outlying cases.

**get**
Context Sentences: 11

if something that is always approved why does it require prior authorization

**request**
Context Sentences: 10

eliminate any regulatory policies requiring preauthorization for patients to receive services related to opioid use disorder.

**use**
Context Sentences: 9

in response to complaints we received from patients about medicare advantage plans in florida, we found at least 6 plans that require prior authorization every month an injection is given.

**seek**
Context Sentences: 6

**support**

**prior authorization**
Score: 289.87
Subthemes: 15

**prior authorization process**

# OPTIMIZING NLP VISUALIZATIONS – ATTEMPT 3

## Themes and Sub-themes with Sentiment (Theme-driven View)

**Human-Centered Design**
Center of Excellence

CCSQ WORLD USABILITY DAY

# OPTIMIZING NLP VISUALIZATIONS – ATTEMPT 4 SUCCESSFUL

## Themes and Sub-themes (Interactive, Theme-driven View)
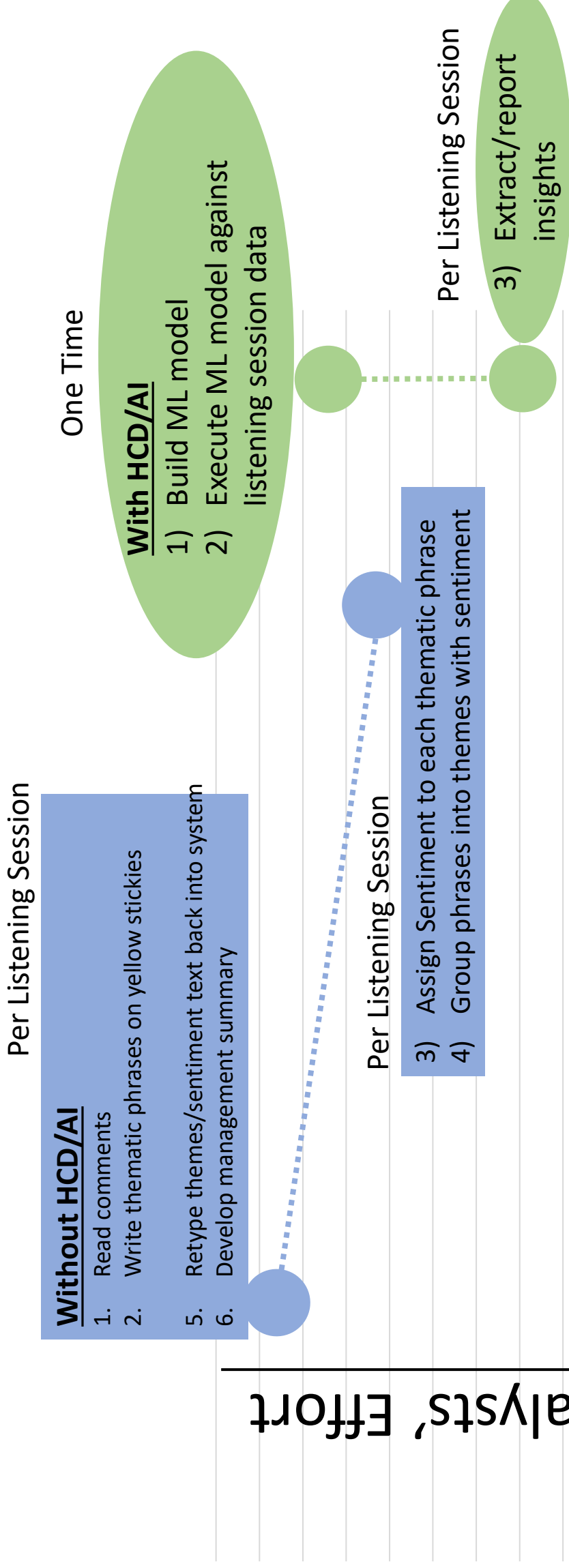


Theme

Sub-themes

Sentence

Interactive Filters

# ATTAINED PROJECT GOALS AND OBJECTIVES

**Goals achieved:**

- Identified a visualization technique that made it possible to interpret listening session data with a large number of themes/sub-themes and their associated sentiment by using an interactive dashboard.

**Human-Centered Design**
Center of Excellence

# HCD Project Outcomes

**Human-Centered Design**
Center of Excellence

# HCD/AI MOVES ANALYSTS FROM LOW-VALUE TO HIGH-VALUE WORK

One Time

**With HCD/AI**
1) Build ML model
2) Execute ML model against listening session data

Per Listening Session
3) Extract/report insights

Per Listening Session

**Without HCD/AI**
1. Read comments
2. Write thematic phrases on yellow stickies
5. Retype themes/sentiment text back into system
6. Develop management summary

Per Listening Session
3) Assign Sentiment to each thematic phrase
4) Group phrases into themes with sentiment

Analysts' Effort

High Value Work

Low Value Work

# HCD HELPS TO FIND GOLD IN UNSTRUCTURED DATA

➢ Artificial Intelligence automates previously manual tasks which:

- Achieve high reliability (>80%) accuracy
- Minimize human bias
- Accelerate the speed of review for each listening session
- Move analysts from low-value to high-value work

➢ Interactive visualizations are far superior to static visualizations by maximizing human-centered design to target insights, including:

- Numerous filters
- Drill-down capability (theme->sub-theme->sentiment)

➢ Accurate, SME-driven labelling of machine learning training data is crucial to modeling success

**Human-Centered Design**
Center of Excellence

CCSQ WORLD USABILITY DAY

# ADDITIONAL HCD CONSIDERATIONS/PUZZLES

## Themes/Sentiment

- Static Data Set vs Changing Data Set
  - Changing data sets increase complexity and tracking overhead
- Subject Matter Expert Data Labelling

## Visualizations

- Multiple Months of Listening Session Data
  - Increases visualization complexity
  - How to track the review process as data changes over months

**Tantus Technologies Inc.**

**Chris Schilstra, MBA, SAFe, CSPO, CSM**

Innovation Program Manager, PM3

**Mobile:  410-440-9003**
**Office:  410.907.8200  x 131**

**cschilstra@tantustech.com**
**www. tantustech.com**

CCSQ WORLD USABILITY DAY

**Human-Centered Design** Center of Excellence

# THANK YOU

Tantus Technologies, Inc.