



**Prodapt** powering  
global telecom

## Making of an intelligent Virtual Agent to transform Customer Experience

Improves precision, recall & accuracy of NLU model

Credits

Sathya Ramana Varri C

Prashanth Suresh Babu

Sarvagya Nayak

# Virtual agents not able to stand up to consumers' expectations

51% of US consumers don't have faith in VA's ability to respond correctly

Chatbots and virtual agents (VA) have high expectations in terms of customer engagements and overall customer experience. That's why [Business Insider](#) claims that **by 2020, 80% of the organizations will be using virtual agents.**

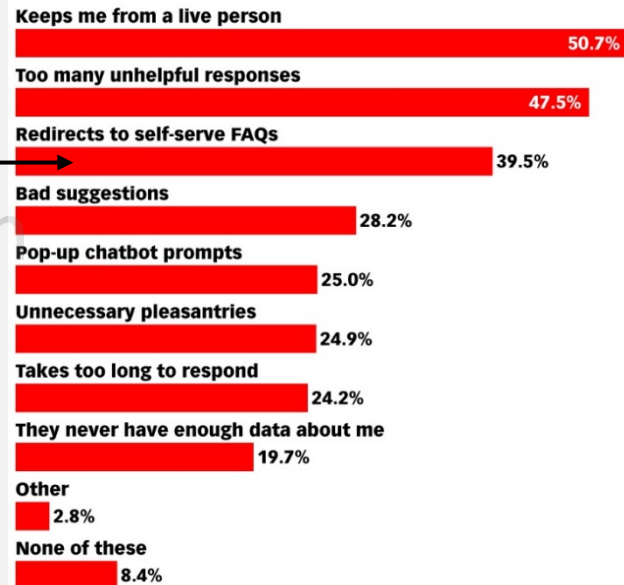
But the end consumers don't have faith in them. In fact, **51%** of the US population thinks that VAs are a hindrance that keeps them from connecting to a live agent. 41% of respondents feel VAs don't provide enough detailed solutions and 37% feel they are generally not helpful.

The **major reason for the failure of these VAs** to satisfy consumers lies in their inability to identify the right intents. This, in turn, is the effect of wrong or inadequate training of VA's natural language understanding (NLU) engine.

Most often the identification of training data is done manually which is not enough. This insight talks about developing a **machine-learning (ML)-based Intent Analyzer tool**, which can identify the most effective data set for NLU training.

## Challenges of Using Chatbots According to US Internet Users, May 2018

% of respondents



Note: ages 18+  
Source: Helpshift, May 31, 2018

238508

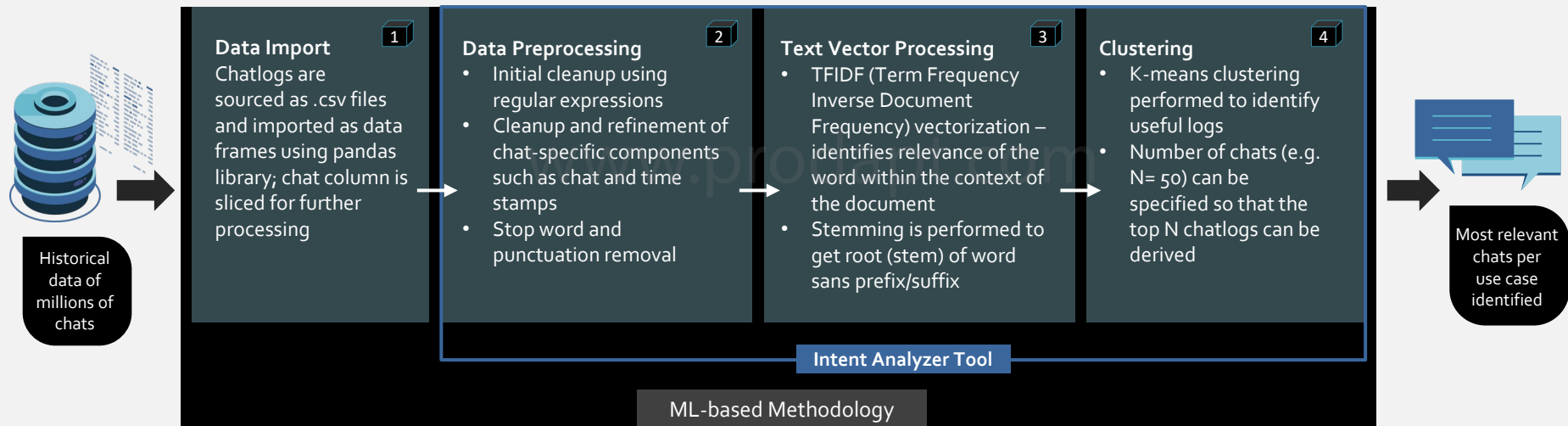
www.eMarketer.com

# Machine learning-based Intent Analyzer tool identifies most relevant representative examples for NLU training

The conventional approach of identifying training data for VA NLU depends heavily on DSP's internal process experts. It involves choosing the most relevant few hundred examples of millions of historical chat. But, it is crippled with inefficiencies because it:

- Lacks coverage of all the examples needed for training
- Makes way for manual biases
- Highly time-consuming

Developing a ML-based intent analyzer tool is the most optimal approach for identifying representative training examples. This tool should perform:



The subsequent slides give details on performing the above 4 steps in the best possible manner.

# Data Import - Templatize the input to optimize the most time-consuming step

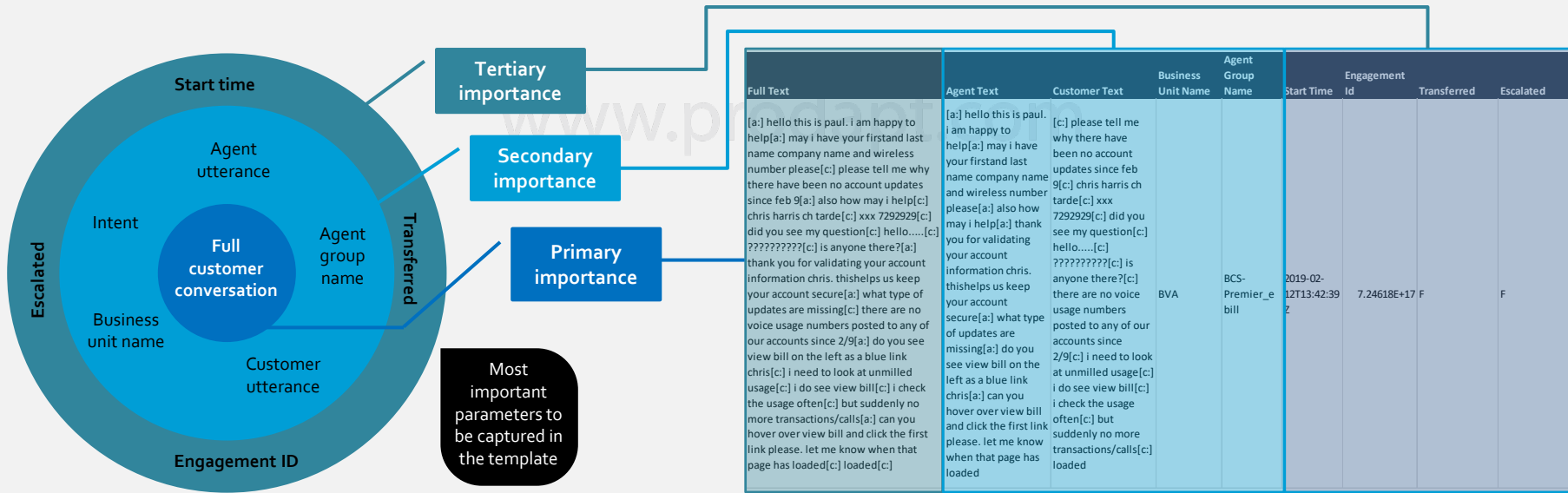
- 1
- 2
- 3
- 4

This step involves sourcing historical data from chatlogs for respective intents/use cases. Millions of chats are picked and imported to identify the most relevant handful of them.

## Recommendations

### Templatize

It is extremely important for the input data to be in a standard format. To ensure this standardization, it is recommended to create a **template** for it



Full Text	Agent Text	Customer Text	Business Unit Name	Agent Group Name	Engagement Start Time	Engagement Id	Transferred	Escalated
[a:] hello this is paul. i am happy to help[a:] may i have your first and last name company name and wireless number please[c:] please tell me why there have been no account updates since feb 9[a:] also how may i help[c:] chris harris ch tarde[c:] xxx 7292929[c:] did you see my question[c:] hello.....[c:] ??????????[c:] is anyone there?[a:] thank you for validating your account information chris. this helps us keep your account secure[a:] what type of updates are missing[c:] there are no voice usage numbers posted to any of our accounts since 2/9[a:] do you see view bill on the left as a blue link chris[c:] i need to look at unmlled usage[c:] i do see view bill[c:] i check the usage often[c:] but suddenly no more transactions/calls[a:] can you hover over view bill and click the first link please. let me know when that page has loaded[c:] loaded[c:]	[a:] hello this is paul. i am happy to help[a:] may i have your first and last name company name and wireless number please[a:] also how may i help[a:] thank you for validating your account information chris. this helps us keep your account secure[a:] what type of updates are missing[a:] do you see view bill on the left as a blue link chris[a:] can you hover over view bill and click the first link please. let me know when that page has loaded	[c:] please tell me why there have been no account updates since feb 9[c:] chris harris ch tarde[c:] xxx 7292929[c:] did you see my question[c:] hello.....[c:] ??????????[c:] is anyone there?[c:] there are no voice usage numbers posted to any of our accounts since 2/9[c:] i need to look at unmlled usage[c:] i do see view bill[c:] i check the usage often[c:] but suddenly no more transactions/calls[c:] loaded	BVA	BCS-Premier_e bill	2019-02-12T13:42:39Z	7.24618E+17	F	F

Sample from a template

# Data Import – Remove random noise and flatten the input file to remove metadata

- 1
- 2
- 3
- 4

## Recommendations

### Remove random noise/white noise

Seasonality in data can lead to wrong inferences. For example, higher call drops or lower speeds during Thanksgiving or Christmas. To reduce that impact, choose the data set spread over a larger time period like 9-12 months.

### Key-value pair

For efficient separation of metadata, flatten the file into excel file or other simpler formats

```
<?xml version="1.0" encoding="ISO-8859-1"?>
<autnresponse xmlns:autn="http://schemas.autonomy.com/act/">
  <action>QUERY</action>
  <response>SUCCESS</response>
  <responsedata>
    <autn:numhits>100</autn:numhits>
    <autn:hit>
      <autn:reference>CONSUMER...226088856.319658996</autn:reference>
      <autn:id>13966180</autn:id>
      <autn:section>0</autn:section>
      <autn:weight>96.00</autn:weight>
      <autn:database>Explore</autn:database>
      <autn:title>ONSUSER...226088856.319658996</autn:title>
      <autn:content>
        <RECORDING RECORDING_ID="319658996" RECORDING_ID="226088856" RECORDING_SOURCE="NSUSER" DB_NAME="Explore"
          UNIQUE_REFERENCE_ID="ONSUSER...226088856.319658996" REFERENCE_ID="ONSUSER...226088856.319658996">
          <TTT>
            <AGENT_NAME="Chapman, Lasonya (BIR)-Jc3534" AGENT_ID="43392"/>
            <INSTALK_PERCENTAGE="0" CROSSTALK_DURATION="0" LOW_EMOTION_DURATION="171130" MEDIUM_EMOTION_DURATION="0"
              TION_DURATION="0" DNIS="4981013,18003310500,8003310500" ANI="2083213600,4311032" GROUP="" SILENCE="93" NUM_HOLD="1"
              WHOLE="265" DURATION="265" DATE_TIME="2019-04-23 23:31:10 GMT"/>
          </TTT>
        </autn:content>
      </autn:hit>
    </autn:responsedata>
  </response>
</autnresponse>
```

Irrelevant metadata deeply embedded in source chat file

autn:cont	autn:cont	autn:cont	autn:cont	autn:cont	autn:cont	autn:cont	autn:cont	autn:cont	autn:cont	autn:cont	autn:cont	autn:cont	autn:cont
1.56E+09	["Henders	2019-04-2	1.56E+09	1	-100	2.26E+08	["i'm sorry	-8	["thanks f	["0.9', '1.1'	t and t a l		
1.56E+09	["Moh	2019-04-2	1.56E+09	1	-100	2.26E+08		-8	["you are r	["1', '0.92',	are you ar		
1.56E+09	2019-04-2	1.56E+09	1	-100	2.26E+08			-8	["okay mh	["0.8', '0.8'	them a da		
1.56E+09	2019-04-2	1.56E+09	1	-100	2.26E+08		["oh no', 'r	-8	["i'm doin	["1.02', '0.8'	hi my nar		
1.56E+09	2019-04-2	1.56E+09	1	-100	2.26E+08		["okay it',	["0.8', '1.3'	the my na				
1.56E+09	2019-04-2	1.56E+09	1	-100	2.26E+08		["i'm sorry	-8	["thanks f	["0.9', '0.64'	thanks for		
1.56E+09	2019-04-2	1.56E+09	1	-100	2.26E+08		["i doubt',	-8	["an issue	["0.65', '0.8'	the high i'		
1.56E+09	2019-04-2	1.56E+09	1	-100	2.26E+08		["okay no	-8	["wanna g	["1', '1', '1',	you come		
1.56E+09	2019-04-2	1.56E+09	1	-100	2.26E+08		calling my	-8	["verse ok	["0.96', '0.5'	thank you		
1.56E+09	2019-04-2	1.56E+09	1	-100	2.26E+08		i'm sorry	-8	i'm not a	["0.864', '1'	the reaso		
1.56E+09	2019-04-2	1.56E+09	1	-100	2.26E+08		["oh no', "	-8	["technica	["0.8', '1.1'	thank alla		
1.56E+09	2019-04-2	1.56E+09	1	-100	2.26E+08		["imperso	-8	["that wou	["0.672', '0'	the them		
1.56E+09	2019-04-2	1.56E+09	1	-100	2.26E+08		["wrong al	-8	["support	["0.9', '0.65'	thank you		
1.56E+09	2019-04-2	1.56E+09	1	-100	2.26E+08		["i'm havi	-8	["that's fr	["0.96', '0.5'	the i'm ye		
1.56E+09	2019-04-2	1.56E+09	1	-100	2.26E+08		["trouble t	["-2', '-8'	["can help	["0.65', '0.8'	hey this is		
1.56E+09	2019-04-2	1.56E+09	1	-100	2.26E+08			-8	["this ipad	["-2', '-8'	["home ph	["1.02', '1',	the yes w
1.56E+09	2019-04-2	1.56E+09	1	-100	2.26E+08		i'm sorry l	-8	["the num	["1', '0.8', '1'	this is trie		
1.56E+09	2019-04-2	1.56E+09	1	-100	2.26E+08		lost remi	-8	["wanna si	["0.8', '1.06'	the hi um		

Metadata and other less useful data segregated in the flattened file

### Reducing import time

By performing parallel processing and avoiding overloading memory

# Data Preprocessing - Leverage raw text preprocessing, regular expression and lemmatization



- 1
- 2
- 3
- 4

The step involves initial clean up of the chat data by removing chat specific components such as timestamps, stop-words and punctuations.

## Recommendations

### Special character processing and text analysis

For removing special characters and analyzing text at high level perform 'raw text preprocessing' and 'regular expressions'

### Raw text preprocessing

Removes chat specific notations and special characters which doesn't add any value to analysis. E.g – timestamps, special characters, etc.

### Regular expression

Segregates numerals from alphabets and retains only special strings of alphanumeric values

[c:] i wish to make a payment arrangement for 13.26 on january 25th xxxx[c:] i have a question about a payment arrangement[c:] ##url#https://www.abcde.com/esupport/article.html#!/iptv/km1025834

i wish to make a payment arrangement for 13.26 on january 25th xxxx i have a question about a payment arrangement url https://www.abcde.com/esupport/article.html#!/iptv/km1025834

my bill was 260\$ and some change will my new bill be around 280\$? [c:] i always thought it was lte since my phone has a symbol of 4g lte

my bill was xxxx\$ and some change will my new bill be around xxxx\$? i always thought it was lte since my phone has a symbol of 4g lte

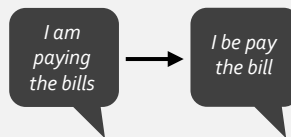
### Lemmatization or stemming

Focuses on reducing the words to their root words

### Lemmatization

the stemmed version of "I am paying the bills" is

Word	Root
Paying, pays, paid	Pay
Am, are, is	Be
Bill, bills, bill's	Bill



### Rare word removal

They create noise by their association with other words. They might not be rare, but their usage in certain context can be misleading. E.g. – erroring, revert, captive, hill, angel.

# Text Vector Pre-processing - Leverage term frequency inverse document frequency (TFIDF) for most effective vectorization

Text vector preprocessing helps in understanding the importance of words as per the relative context. It focusses on the difference in relevance in different circumstances.

## Recommendations

### Choice of vectorization method

Choose a method based on the type of text data. It is recommended to choose 'term frequency-inverse document frequency (TFIDF) vectorization' since it considers the relative importance of a word in each context.

### Term Frequency-Inverse Document Frequency (TFIDF)

- TFIDF ensures that the chats are selected according to their relative importance
- More than their overall significance in everyday usage, it measures how critical they are in the context of the chat log corpus being analyzed
- This helps in identifying the most relevant chats as per the intent

#### N-grams and other multiword usage

Certain words have an entirely different meanings when used in combination with a few other words. Such words should be configured appropriately.

*E.g. – the words "payment" and "arrangement" have different sense and relevance when used individually. But upon using collectively as "payment arrangement" it conveys some other meaning.*

#### Hyper-parameter tuning

Tunes the parameters of the vectorization algorithm to optimize the output.

- E.g.*
- *Max\_df* – sets the acceptable upper limit of frequency.
    - *E.g. – 'IPTV max\_df = .85' - Chat containing "IPTV" more than 85% of time will be ignored.*
  - *Min\_df* – sets the acceptable lower limit of frequency
    - *E.g. – 'IPTV min\_df = .20' - Chat containing "IPTV" less than 20% of time will be ignored*
  - *Max\_features* – defines the vocabulary size
    - *E.g. – 'Max\_features = 10,000' - This limits the number of words in vocabulary to 10,000 words*

### TFIDF Input

E.g. – How to close my account,  
Close my account,  
Delete my account,  
Close account



### TFIDF Output

(0, 157)	0.6470729031869088
(0, 6)	0.5314083612599048
(0, 213)	0.4035924769447586
(0, 447)	0.3688020120578553

# Clustering - Perform k-means clustering technique for effective classification

Clustering ensures that the top-N chats (where N is variable depending on business/NLU needs) are derived. These can be quickly analyzed to identify utterances, intents and entities. Additional ML processing such as entity or intent recognition can also be performed if required. All this results in significant time and effort saving.

## Recommendations

### Choice of clustering technique

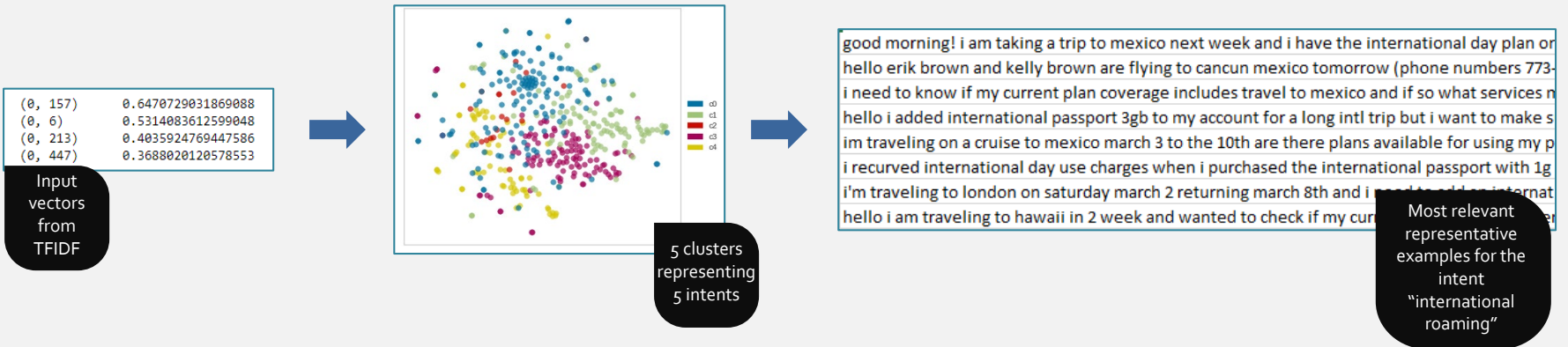
Choose the technique that can work on huge volume of data like millions of customer chats. It is recommended to use k-means clustering for such a volume.

### One-on-one mapping

Ensure one chat is mapped with only one intent i.e. avoid overlapping

### Intent-specific scaling

Ability to scale the number of top-N use-cases based on intent call-volume (by varying the number of clusters): this enables the number of representative samples to be adjusted based on whether a given intent has more or less volume





# Key takeaways

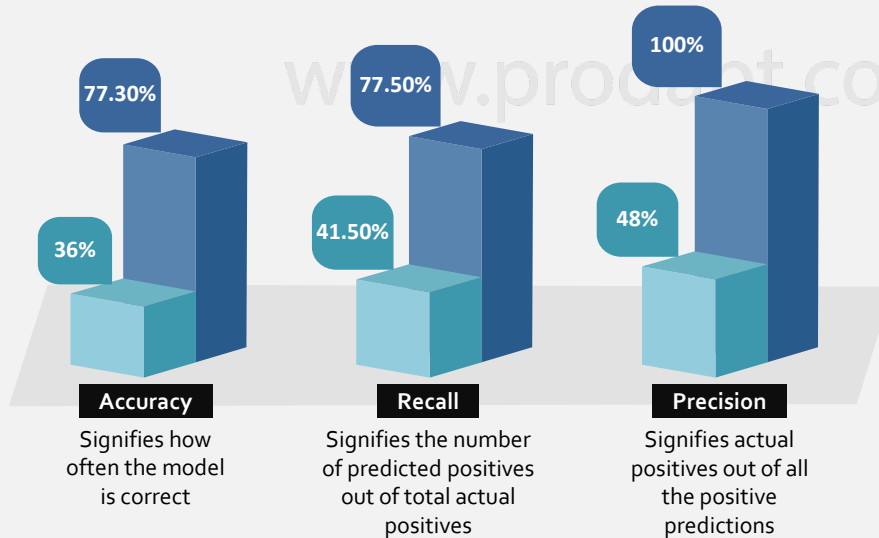
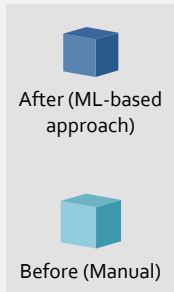
Leveraging a machine learning-based approach for identification of training data for NLU can have following benefits



## NLU confidence

The number of use cases crossing confidence threshold can increase by **160-180%**.

**Confidence threshold** – the minimum confidence level configured in VA below which it can't process the chat and transfers it to the live agent



## Time efficiency

Can save up to **97%** of time in identifying the examples



## Transfer to live agents

Can reduce by almost **80%**

# Get in touch

## USA

**Prodapt North America**  
Tualatin: 7565 SW Mohawk St.,  
Phone: +1 503 636 3737

**Dallas:** 1333, Corporate Dr., Suite 101, Irving  
Phone: +1 972 201 9009

**New York:** 1 Bridge Street, Irvington  
Phone: +1 646 403 8161

## CANADA

**Prodapt Canada Inc.**  
Vancouver: 777, Hornby Street,  
Suite 600, BC V6Z 1S4  
Phone: +1 503 210 0107

## UK

**Prodapt (UK) Limited**  
Reading: Davidson House,  
The Forbury, RG1 3EU  
Phone: +44 (0) 11 8900 1068

## EUROPE

**Prodapt Solutions Europe &  
Prodapt Consulting BV**  
Rijswijk: De Bruyn Kopsstraat 14  
Phone: +31 (0) 70 4140722

**Prodapt Germany GmbH**  
München: Brienner Straße, 80333  
Phone: +31 (0) 70 4140722

**Prodapt Switzerland GmbH**  
Zürich: Mühlebachstrasse 54,  
8008 Zürich

## SOUTH AFRICA

**Prodapt SA (Pty) Ltd.**  
Johannesburg: No. 3,  
3rd Avenue, Rivonia  
Phone: +27 (0) 11 259 4000

## INDIA

**Prodapt Solutions Pvt. Ltd.**  
Chennai: Prince Infocity II, OMR  
Phone: +91 44 4903 3000

“Chennai One” SEZ, Thoraipakkam  
Phone: +91 44 4230 2300

IIT Madras Research Park II,  
3<sup>rd</sup> floor, Kanagam Road, Taramani  
Phone: +91 44 4903 3020

**Bangalore:** “CareerNet Campus”  
2<sup>nd</sup> floor, No. 53, Devarabisana Halli,  
Phone: +91 80 4655 7008

# THANK YOU!

The success of any Virtual Agent (VA) depends on the training of its Natural Language Understanding (NLU) model prior to configuration. The challenge is providing the right set of representative examples from historical data for this training. Identifying few hundreds of right example out of millions of historical data is a herculean task. What makes it even more daunting is that this task is usually done by digital service providers (DSPs) manually. This not only makes finding the most suitable examples questionable but also extremely time consuming.

This insight talks about developing a Machine Learning (ML) based tool to identify most appropriate and small data set of representative examples for training. These examples covers maximum scope for the respective intent making the training of NLU highly efficient leading to improved precision, recall and accuracy. The ultimate benefit of this is improved customer experience, containment and reduced abandonment. Since this a tool-based approach, it also saves a lot of time in comparison to manually identifying the training examples. Improved training efficiency in the first time also saves time and efforts in the subsequent re-training.