# New Development Property Prediction

## Property Advertisement Description Classification

Domain

Supervisor： Min Li
　　　　　 (Head of Data Science - Data Quality and Information Extraction)
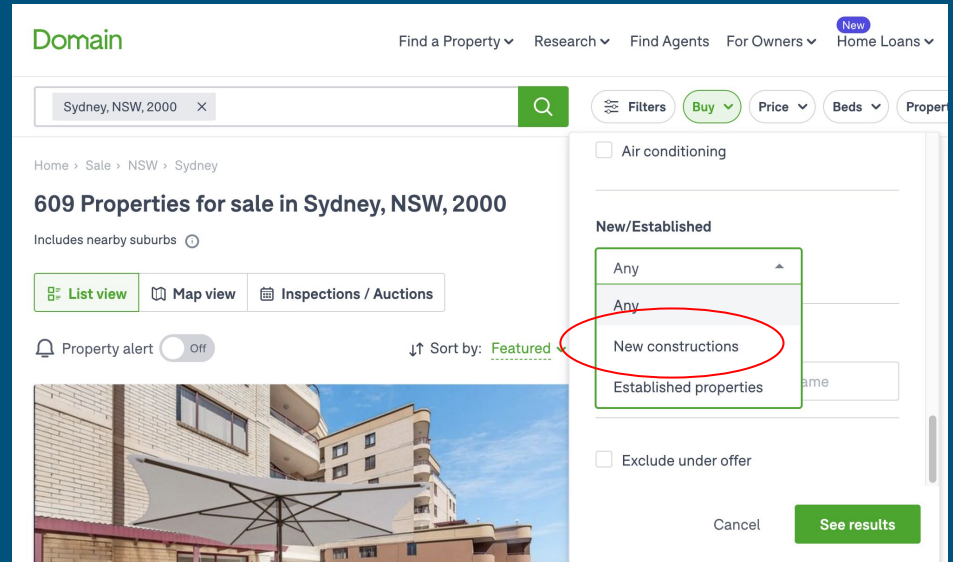Major：　　 Master of Data Science
Name ：　　 Lei Yang

# PROJECT OBJECTIVE

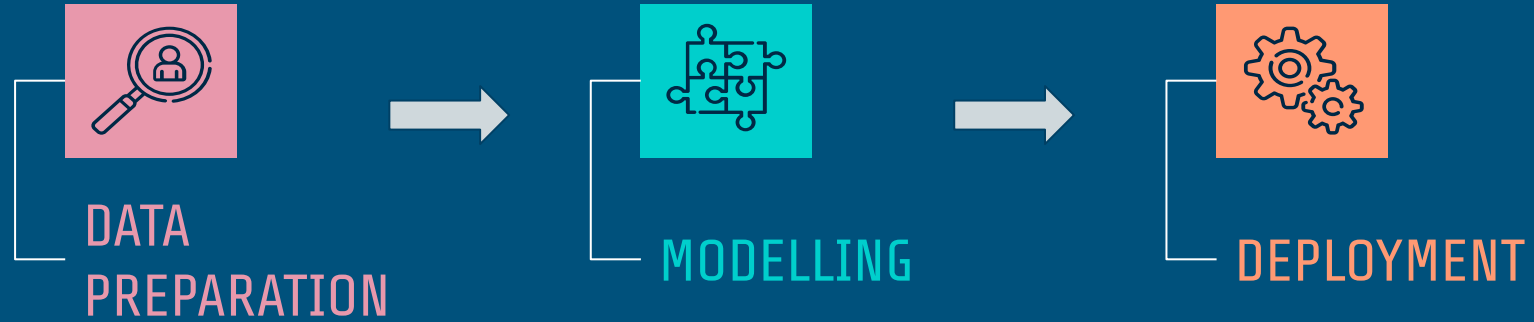**Project Title:** Description Classification - New Development vs. Established

**Project Objective/Goal:** Assist the business in understanding the number of new development properties in the market
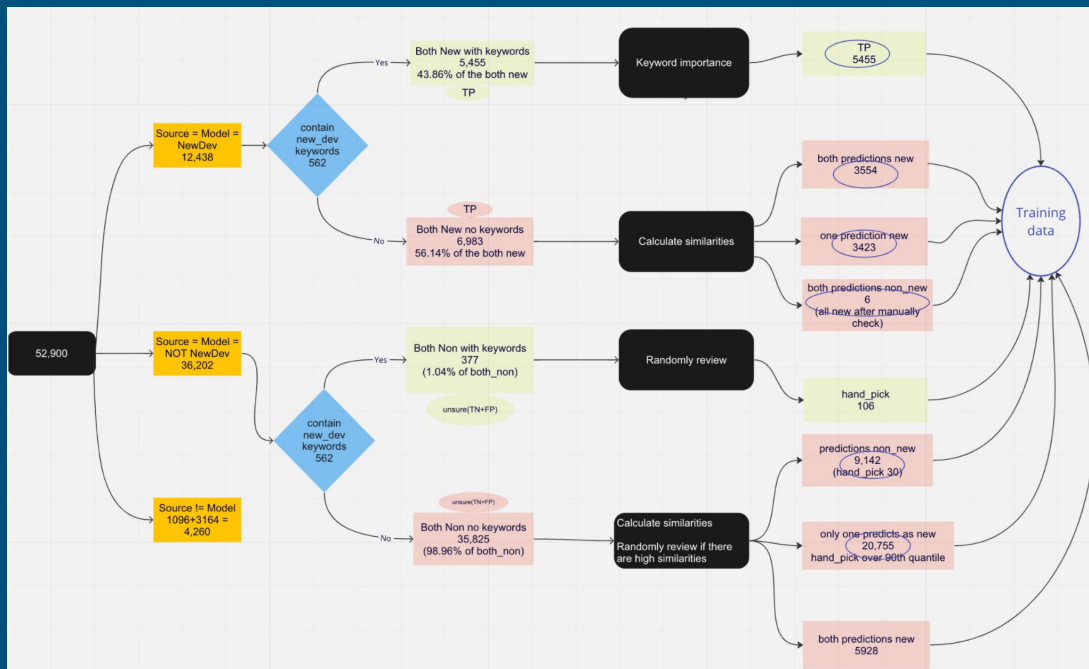
**Problem Description:**

- Improve filtering accuracy
- Improve user experience
- Enhance scalability of profit
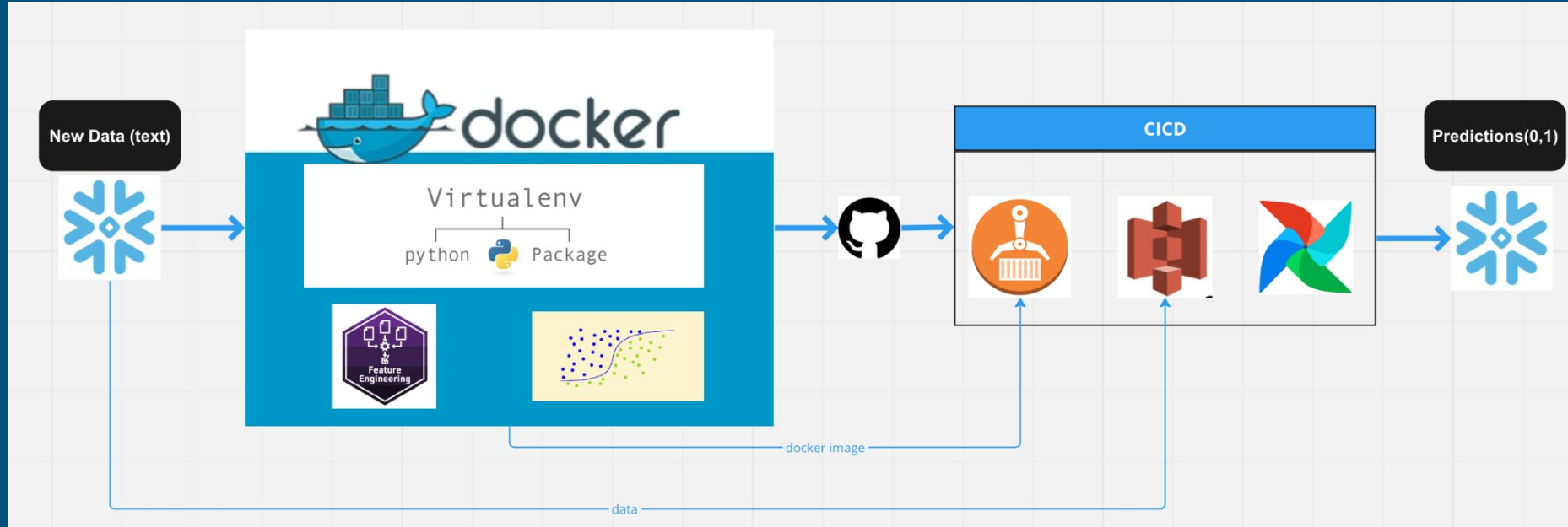
# PROJECT IMPLEMENTATION – DATA PREPARATION

## Flowchart (left)

52,900

Source = Model = NewDev 12,438

Source = Model = NOT NewDev 36,202

Source != Model 1096+3164 = 4,260

contain new_dev keywords 562

- Yes → Both New with keywords 5,455 43.86% of the both new → TP → Keyword importance → TP 5455
- No → Both New no keywords 6,983 56.14% of the both new → TP → Calculate similarities
  - both predictions new 3554
  - one prediction new 3423
  - both predictions non_new 6 (all new after manually check)

contain new_dev keywords 562

- Yes → Both Non with keywords 377 (1.04% of both_non) → unsure(TN+FP) → Randomly review → hand_pick 106
- No → Both Non no keywords 35,825 (98.96% of both_non) → unsure(TN+FP) → Calculate similarities / Randomly review if there are high similarities
  - predictions_non_new 9,142 (hand_pick 30)
  - only one predicts as new 20,755 hand_pick over 90th quantile
  - both predictions new 5928

→ Training data

## Approaches table (top right)

| Approachs | Vectorizer | Algorithm | | Results | Pros & Cons | |
|---|---|---|---|---|---|---|
| Approach I K-means clustering on texts which both have same labels & no keywords data | TFIDF | Cosine similarities | K-means clustering | | Less interpretation | "new" and "non_new" embeddings are in the same group |
| Approach II Comparation of new and non_new similarities | TFIDF | Cosine similarities | Compare new similarities - non_new similarities | | Limitation of TFIDF | Keywords for both new and non_new |
| Approach III zero shot text classification Lbl2Vec | Lbl2Vec | Cosine similarities | Clean outliners | | keywords can only be trained as words not sentences | |
| | | Find centroids | Classification test | | | |

### Approach III detail
```
doc_key                                    1
most_similar_label                       new
highest_similarity_score            0.347093
new                                 0.347093
non_new                             0.096621
index                                    new
LISTING_ID                         1310205166
text         house and land in banksia grove estate great l...
model_label                              new
tokens       house land banksia grove estate great lot land...
prediction                               new
tagged_docs  ([house, and, land, in, banksia, grove, estate...
```

## Text table (bottom right)

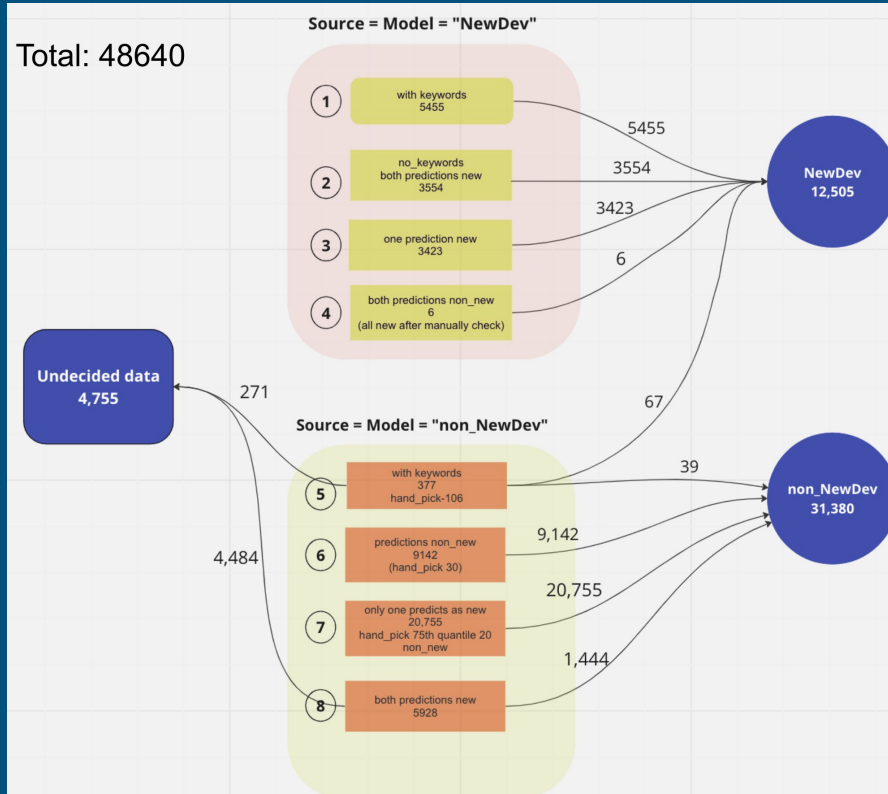| text | source_label | model_label | app2 | app3 |
|---|---|---|---|---|
| pier 32 - luxury harbourside apartments indulge in the ultimate harbourside lifestyle.<br/>34 luxury apartments with unparalleled harbour views. stroll to celebrated restaurants, cafes and boutiques and live the holiday feeling year round.<br/><br/>invest in your future now, to secure your preferred luxury apartment. for further details and conditions, please contact ben pryde on 0410 324 209.<br/><br/>pier 32 offers the unrivalled lifestyle experience with ocean and harbour views from all apartments. ideally located to the best the town has to offer. competitively priced, pier 32 brings luxury and affordability together for the first time in ulladulla.<br/>pier 32 boasts the high quality appliances and finishes available. space & style come together making each 3 bedroom apartment a masterpiece in construction. go to https://pier32ulladulla.com<br/>all apartments enjoy ultra modern open plan living areas adjoining a private entertaining terrace with | non_new | non_new | non_new | new |
| restful parkland views from the spacious apartment in resort style estate a beautifully presented three bedroom apartment with a sunbathed aspect and an open leafy outlook over parklands. boasting resort-style facilities, it is a stone's throw to rhodes shopping centre while offering free-flowing living spaces enhanced by timber flooring and neutral colour palette.<br/><br/>what we like about the property:<br/><br/>+ 3 bedrooms, offering built-in wardrobes and main with ensuite + spacious open plan lounge and dining with air-conditioning and balcony with leafy outlook + beautiful timber flooring | non_new | non_new | non_new | new |

# PROJECT IMPLEMENTATION – MODELLING

# PROJECT IMPLEMENTATION – DEPLOYMENT

# RESULTS/ACHIEVEMENTS – DATA PREPARATION

# RESULTS/ACHIEVEMENTS – MODELLING

# RESULTS/ACHIEVEMENTS – DEPLOYMENT

```
(base) lei.yang@C02FD1D3MD6V ds-nlp-rea-newdev-docker % make app.run
AWS_PROFILE=default \
        .venv/bin/python src/git_test.py sagemaker_code_path_we_ignore --s3-bucket-name decision-science-emr
[nltk_data] Downloading package wordnet to
[nltk_data]     /Users/lei.yang/nltk_data...
[nltk_data]   Package wordnet is already up-to-date!
[nltk_data] Downloading package omw-1.4 to
[nltk_data]     /Users/lei.yang/nltk_data...
[nltk_data]   Package omw-1.4 is already up-to-date!
[nltk_data] Downloading package wordnet to
[nltk_data]     /Users/lei.yang/nltk_data...
[nltk_data]   Package wordnet is already up-to-date!
[nltk_data] Downloading package omw-1.4 to
[nltk_data]     /Users/lei.yang/nltk_data...
[nltk_data]   Package omw-1.4 is already up-to-date!
Complete. Number of text that have been cleaned and tokenized : 10
                                        tokens    0_prob    1_prob  \
0  affordable bedrooms separate lounge kitchen me...  0.995754  0.004246
1  blue sky green tree open space family fun half...  0.000361  0.999639
2  flat lightly tree great beach land quiet stree...  0.980461  0.019539
3  turn key home land package oakdale guarantee f...  0.003831  0.996169
4  touch tradition character fill home perfect ba...  0.996125  0.003875
5  renovate demolish create dream bedroom home of...  0.999110  0.000890
6  floor oceans resort best buy unit currently be...  0.998863  0.001137
7  mimosa home proudly present brook house land p...  0.000548  0.999452
8  bangalow charm truly stun love bangalow home p...  0.997181  0.002819
9  city view lifestyle embrace essential elements...  0.167983  0.832017

   label_pred                  timestamp
0           0 2023-05-24 12:25:20.515265
1           1 2023-05-24 12:25:20.515265
2           0 2023-05-24 12:25:20.515265
3           1 2023-05-24 12:25:20.515265
4           0 2023-05-24 12:25:20.515265
5           0 2023-05-24 12:25:20.515265
6           0 2023-05-24 12:25:20.515265
7           1 2023-05-24 12:25:20.515265
8           0 2023-05-24 12:25:20.515265
9           1 2023-05-24 12:25:20.515265
```
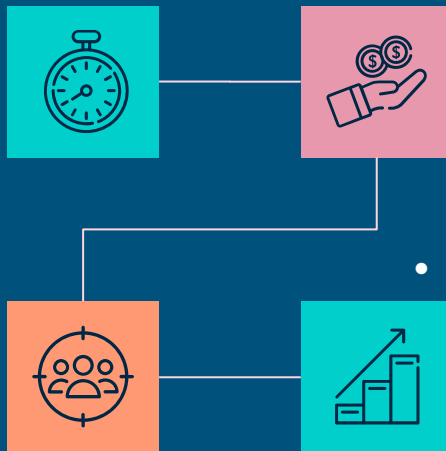
# CONCLUSION

## LESSONS LEARNED

- Modelling is not the only aspect
- Data preparation plays crucial role

## FUTURE IMPROVEMENT

- Regularly update data labels
- Understanding project purpose and overall scope from beginning

## EXPERIENCE

- Real world data science project
- Complete project lifecycle and pipeline

## RECOMMENDATIONS

- The descriptions of new construction properties have maintained a format that attracts potential buyers. This result can assist in refining data cleaning process and model architecture.

- Closely monitor any changes in the description format(more creative or AI-generated). In such cases, leveraging NLP deep learning models could provide an advantage.

THANKS