

# Sentiment Analysis Through the Use of Unsupervised Deep Learning

S7330

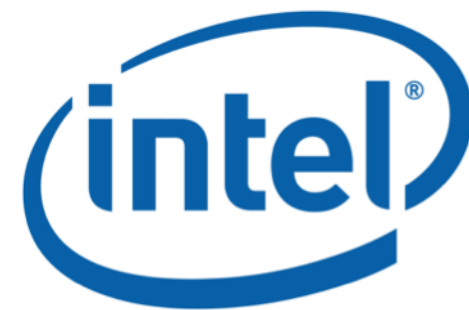
Monday 8<sup>th</sup> May 2017  
GPU Technology Conference

**A. Stephen McGough, Noura Al Moubayed**

Durham University, UK



CUDA™  
RESEARCH  
CENTER



Intel Parallel  
Computing Centre



# Outline

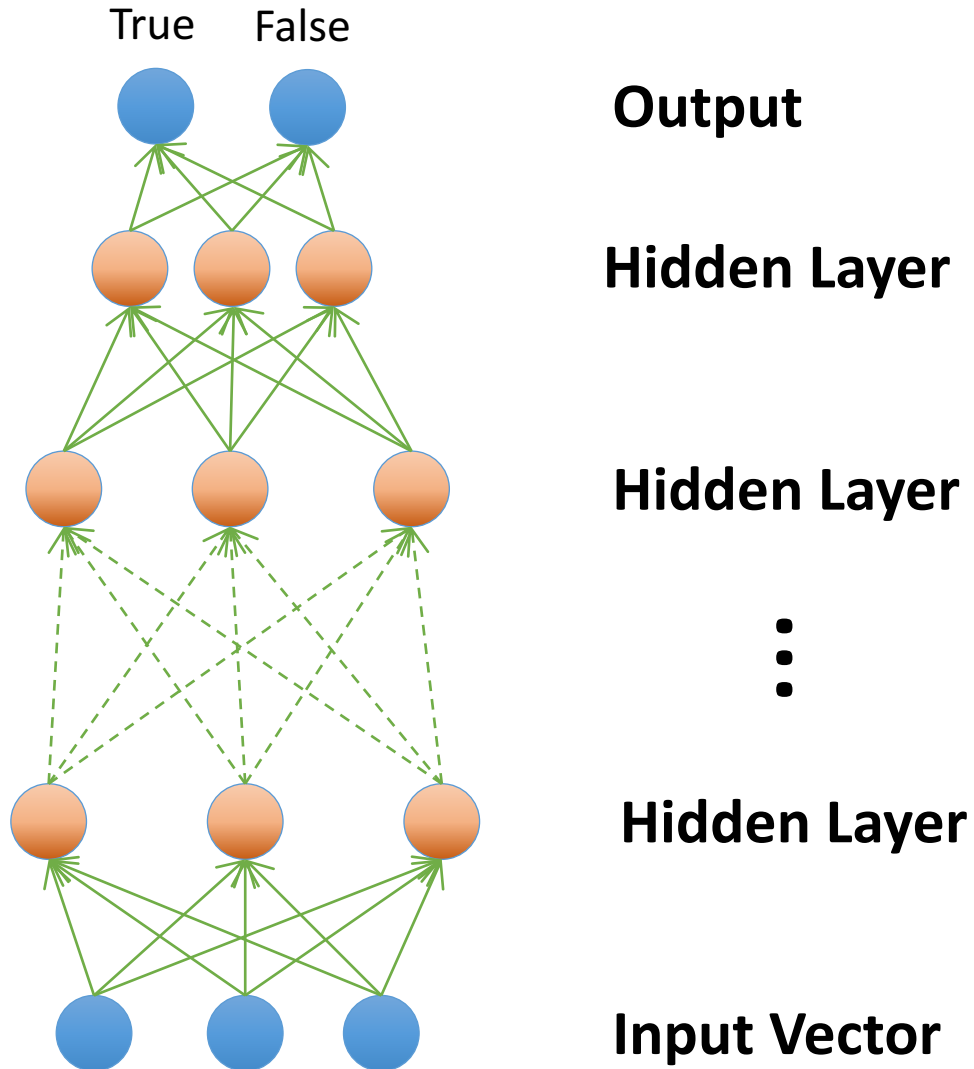
- **The Problem**
- Unsupervised Deep Learning Model
- Text Processing for Topic Modelling
- Detecting Anomalies in Text
- Sentiment Classification



# Outline

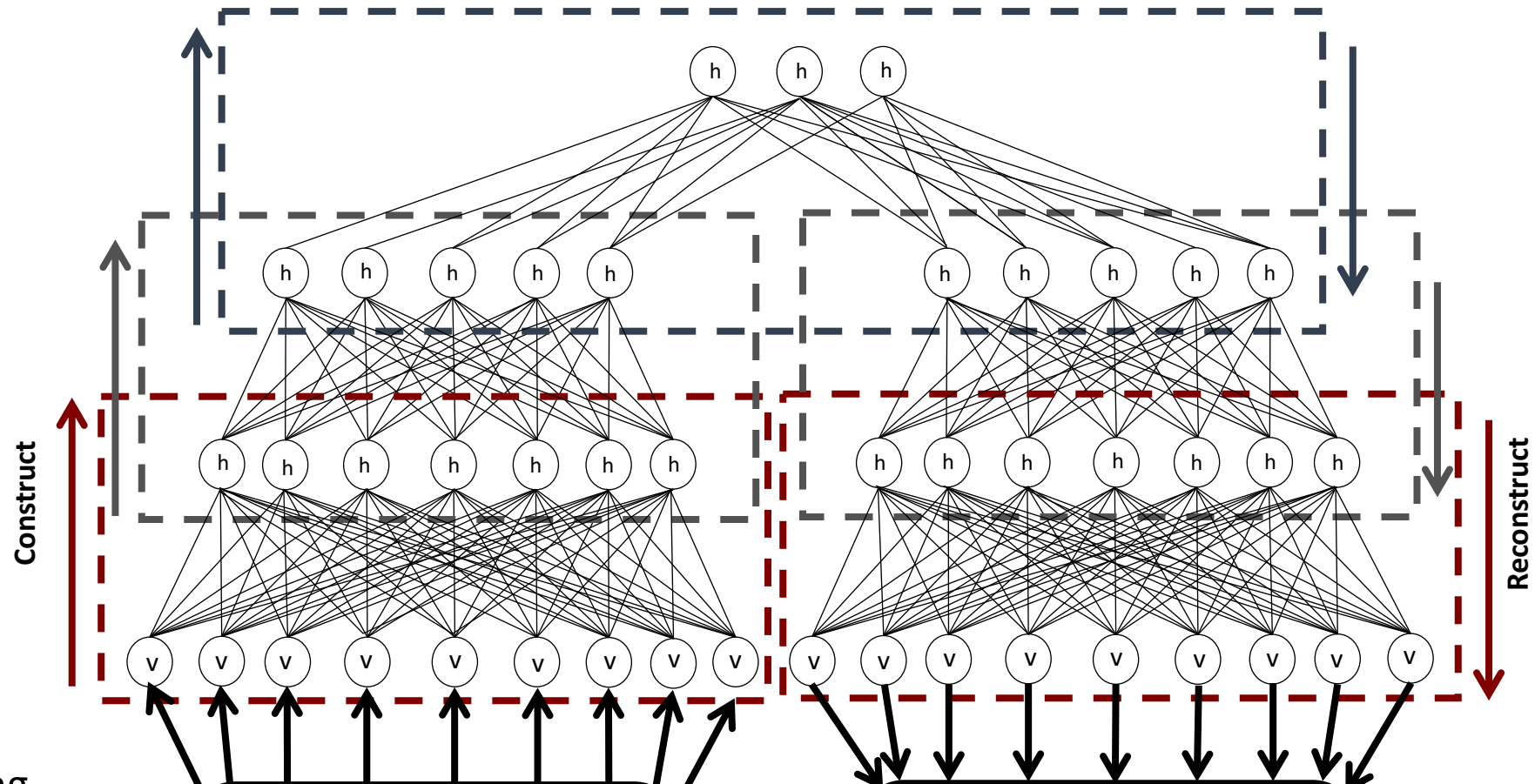
- The Problem
- **Unsupervised Deep Learning Model**
- Text Processing for Topic Modelling
- Detecting Anomalies in Text
- Sentiment Classification

# Deep Learning: Categorization and Regression



But these require labelled data for training

# Anomaly Detection: Unsupervised Deep Learning



Stacked Denoising  
Autoencoder (SDA)

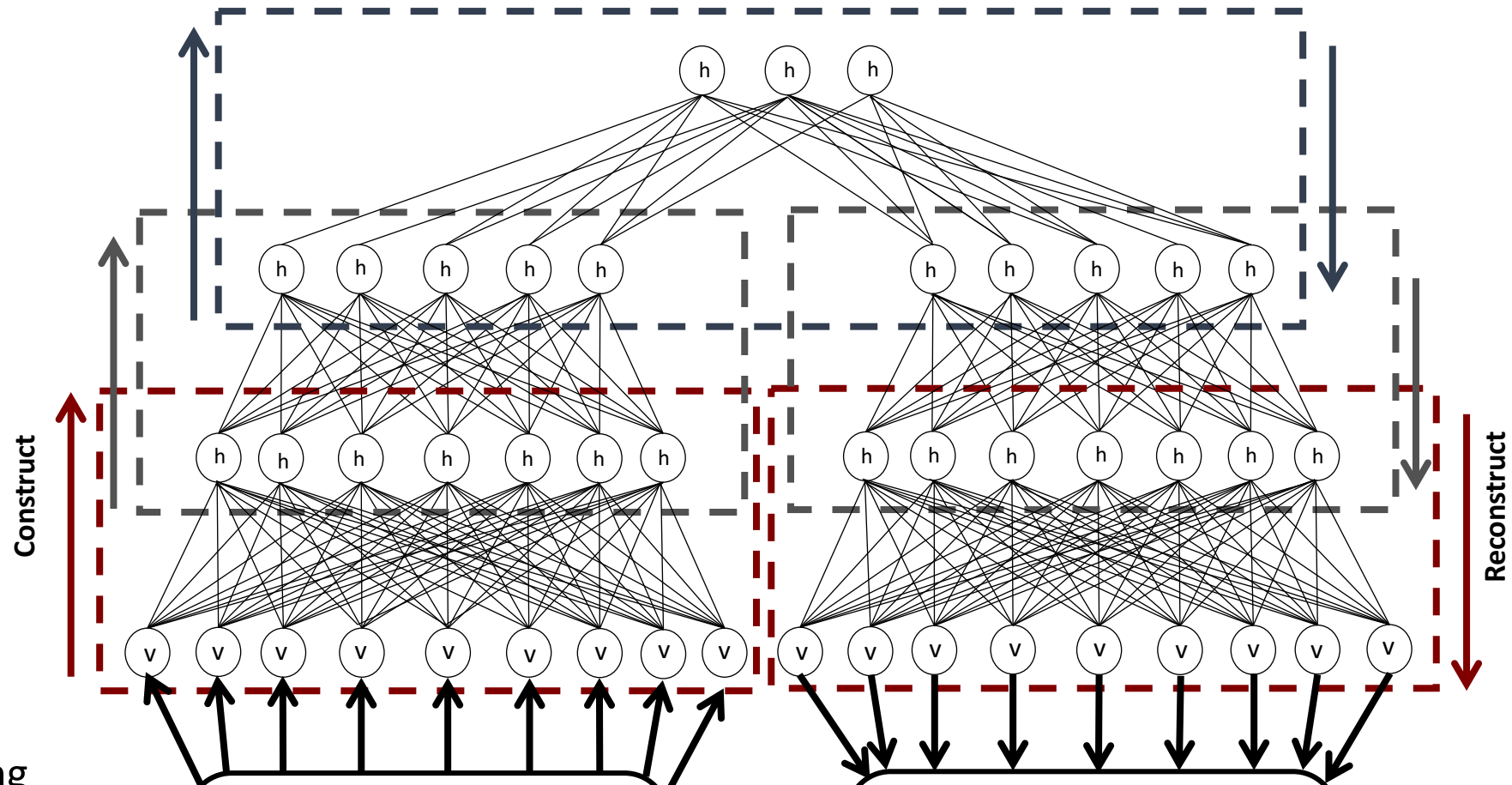
- Trained layer by layer
- Input is corrupted with noise

Input Data  
We're all going to the  
cinema on Saturday

Output Data  
We're all going to the  
theatre on Saturday

Low reconstruction  
error

# Anomaly Detection: Unsupervised Deep Learning



Stacked Denoising  
Autoencoder (SDA)

- Trained layer by layer
- Input is corrupted with noise

Input Data  
Buy Viagra from our  
online meds store

Output Data

By the valley our  
overall medals store

High reconstruction  
error

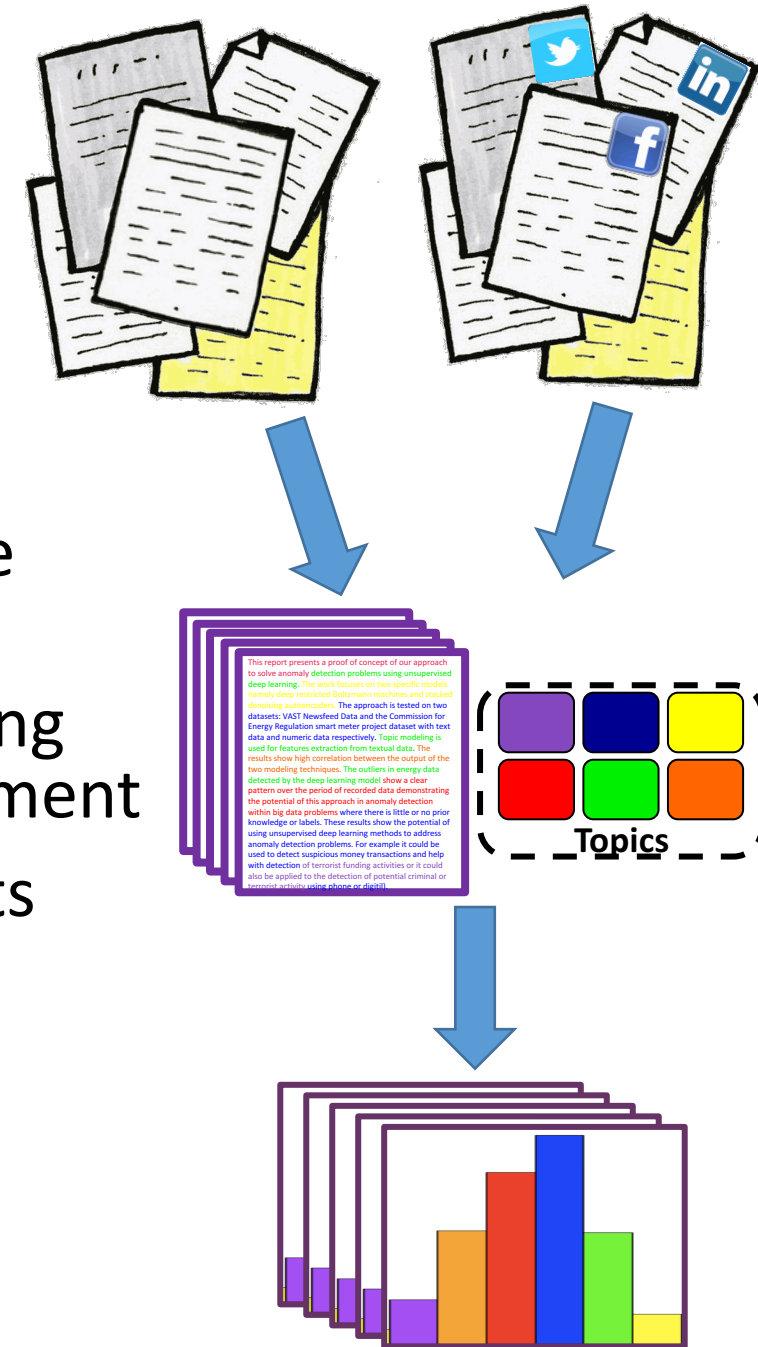


# Outline

- The Problem
- Unsupervised Deep Learning Model
- **Text Processing for Topic Modelling**
- Detecting Anomalies in Text
- Sentiment Classification

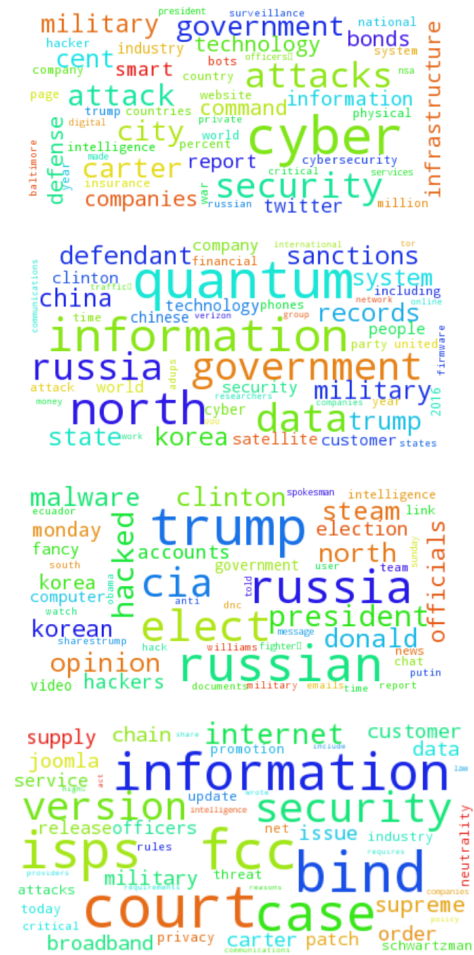
# Probabilistic Topic Modelling

- Unsupervised analysis of text
  - Too many documents to label manually
- Allows us to uncover automatically themes that are latent in a collection of documents
- Same words may have different meanings depending on their co-occurrence with other words in a document
- Statistically identify the topics in a set of documents
  - Which topics appear in which documents
- Statistically identify words in topics
  - Which words co-occur in a document containing topic X
- Document transformed to feature vector of topics

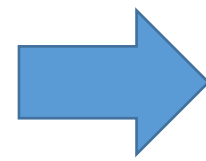
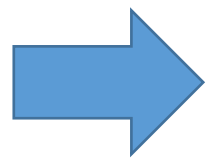


# Topic Modelling: Latent Dirichlet Allocation (LDA)

## Topics



Randomly Select  
topics



Randomly Select  
words from topics

US Govt Data Shows Russia  
Used Outdated Ukrainian PHP  
Malware

The United States government  
earlier this year officially  
accused Russia of interfering  
with the US elections. Earlier  
this year on October 7th, the  
Department of Homeland  
Security and the Office of the  
Director of National Intelligence



# Topic Modelling: Latent Dirichlet Allocation (LDA)

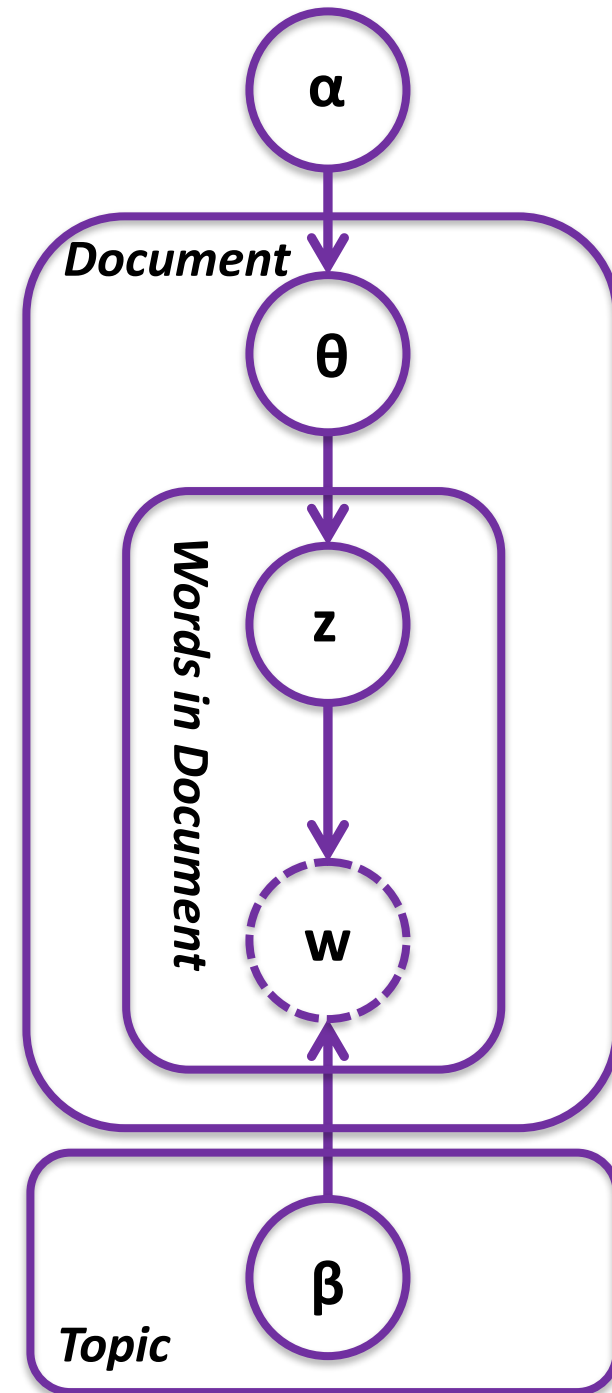
- Each document contains a proportion ( $\theta$ ) of each topic
  - Proportion may be zero
- The probability of each word appearing in a topic ( $\beta$ )
- LDA assumes documents were constructed via:

Words in document =  $N$

Per document topic proportions =  $\theta$  = Dirichlet function ( $\alpha$ )

For each of the  $N$  words  $w_n$ :

- Choose a topic  $z_n \sim \text{Multinomial}(\theta)$
- Choose a word  $w_n$  from topic  $z_n$  (based on probabilities  $\beta$ )
- Need to determine  $\alpha, \beta, \theta$  given we have the documents
  - Solve using variational Bayesian methods or Gibbs sampling



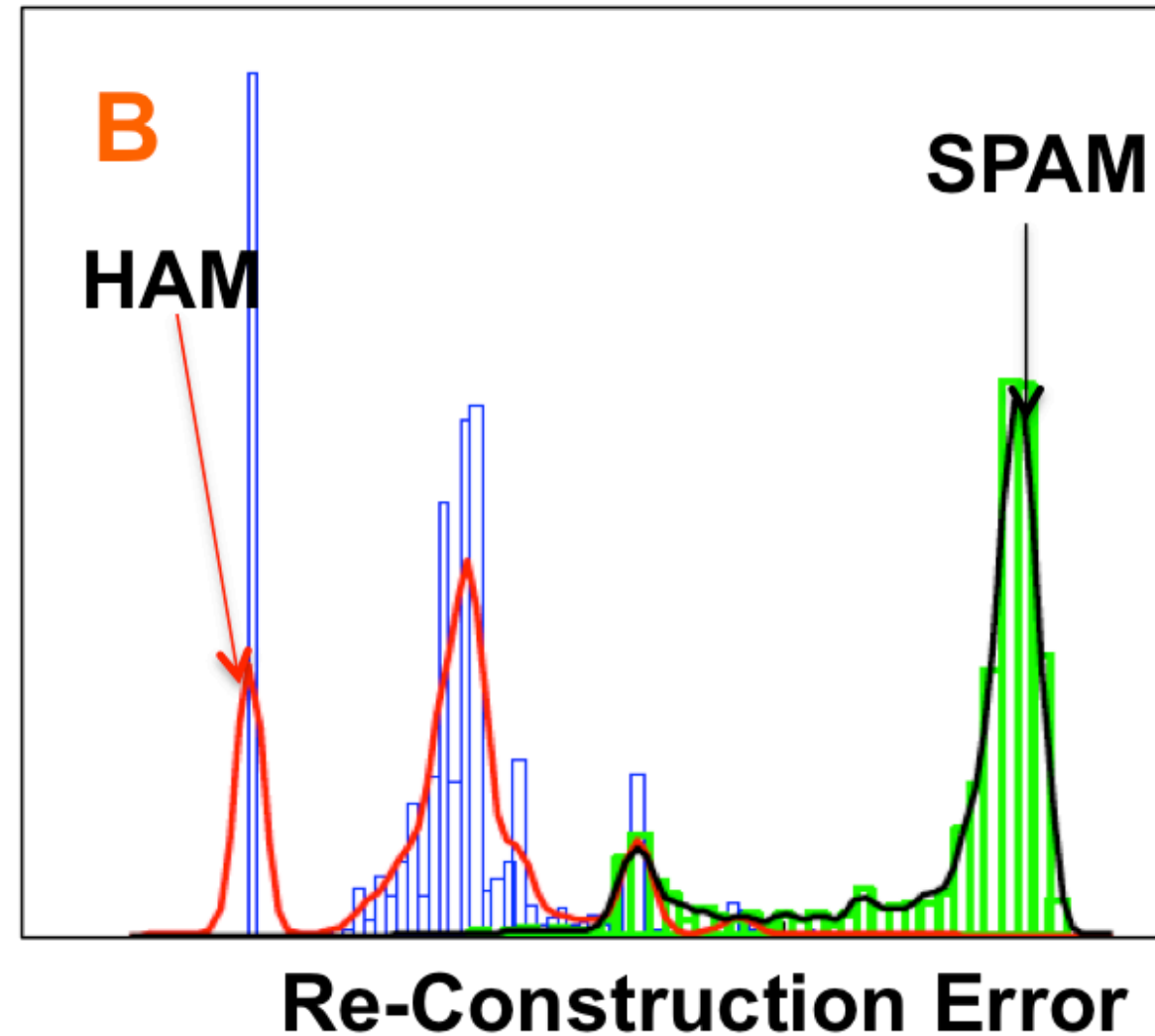
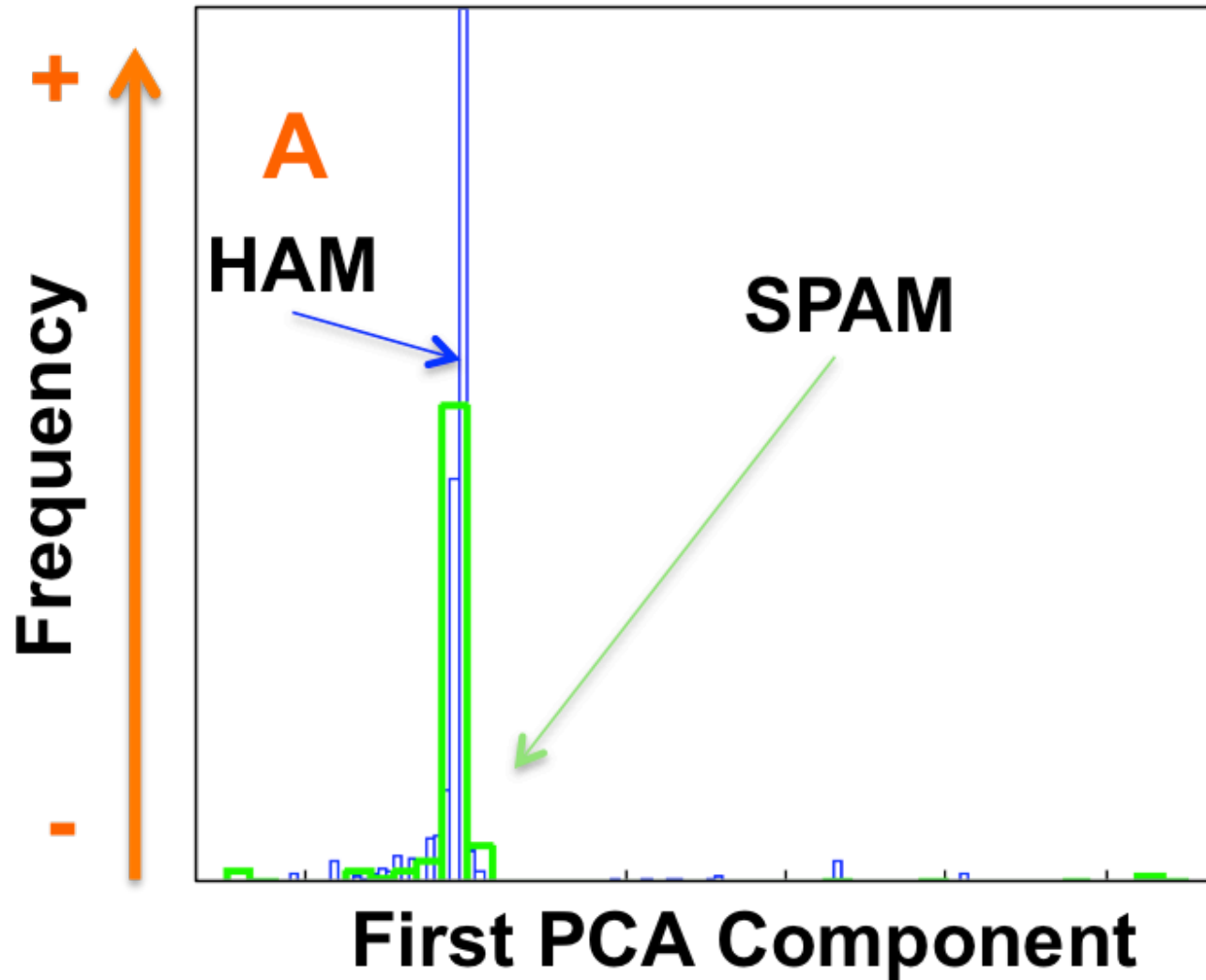
# Outline

- The Problem
- Unsupervised Deep Learning Model
- Text Processing for Topic Modelling
- **Detecting Anomalies in Text**
- Sentiment Classification



# Example: Anomaly Identification

SPAM and HAM in SMS



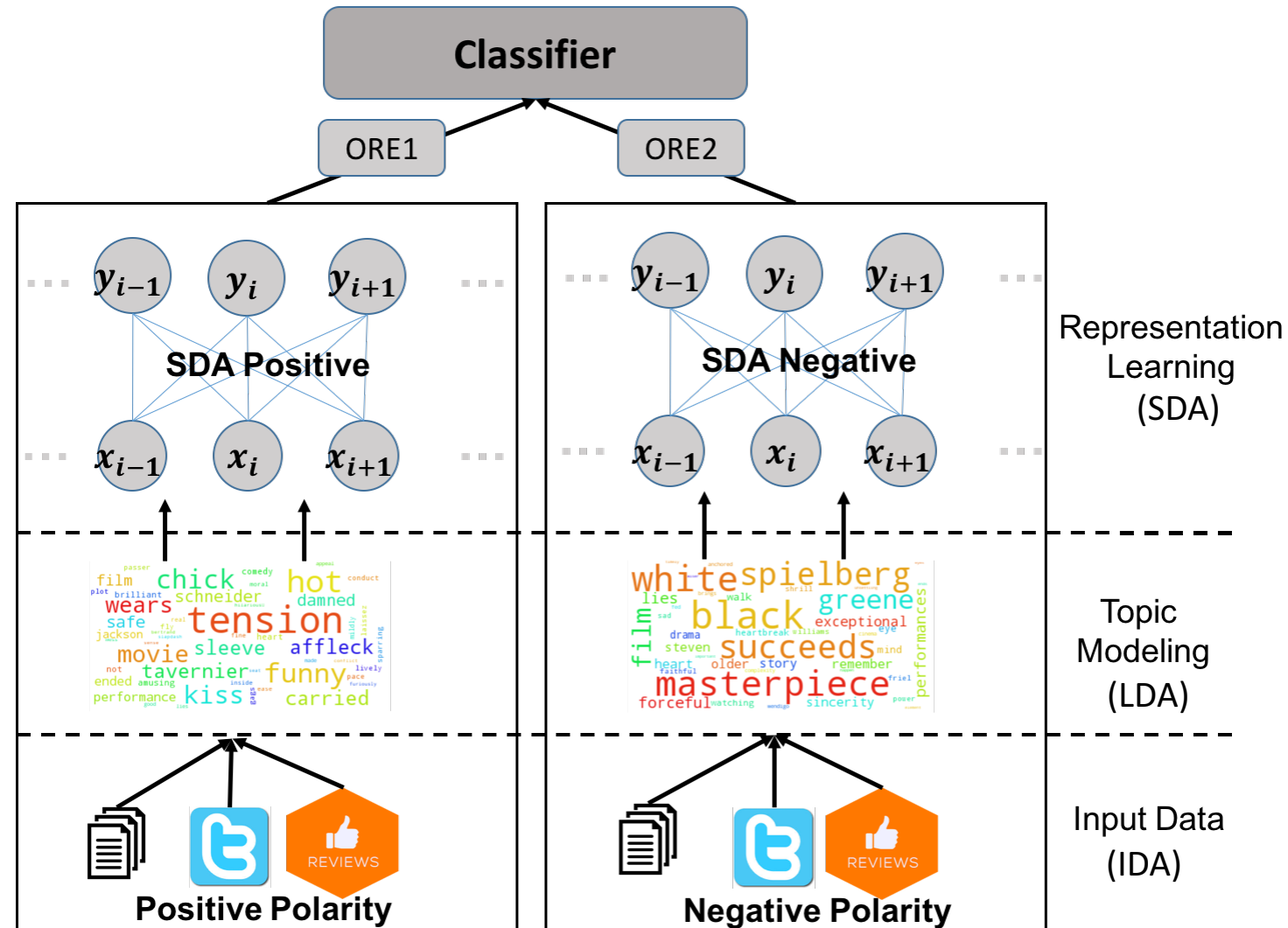
# Outline

- The Problem
- Unsupervised Deep Learning Model
- Text Processing for Topic Modelling
- Detecting Anomalies in Text
- **Sentiment Classification**

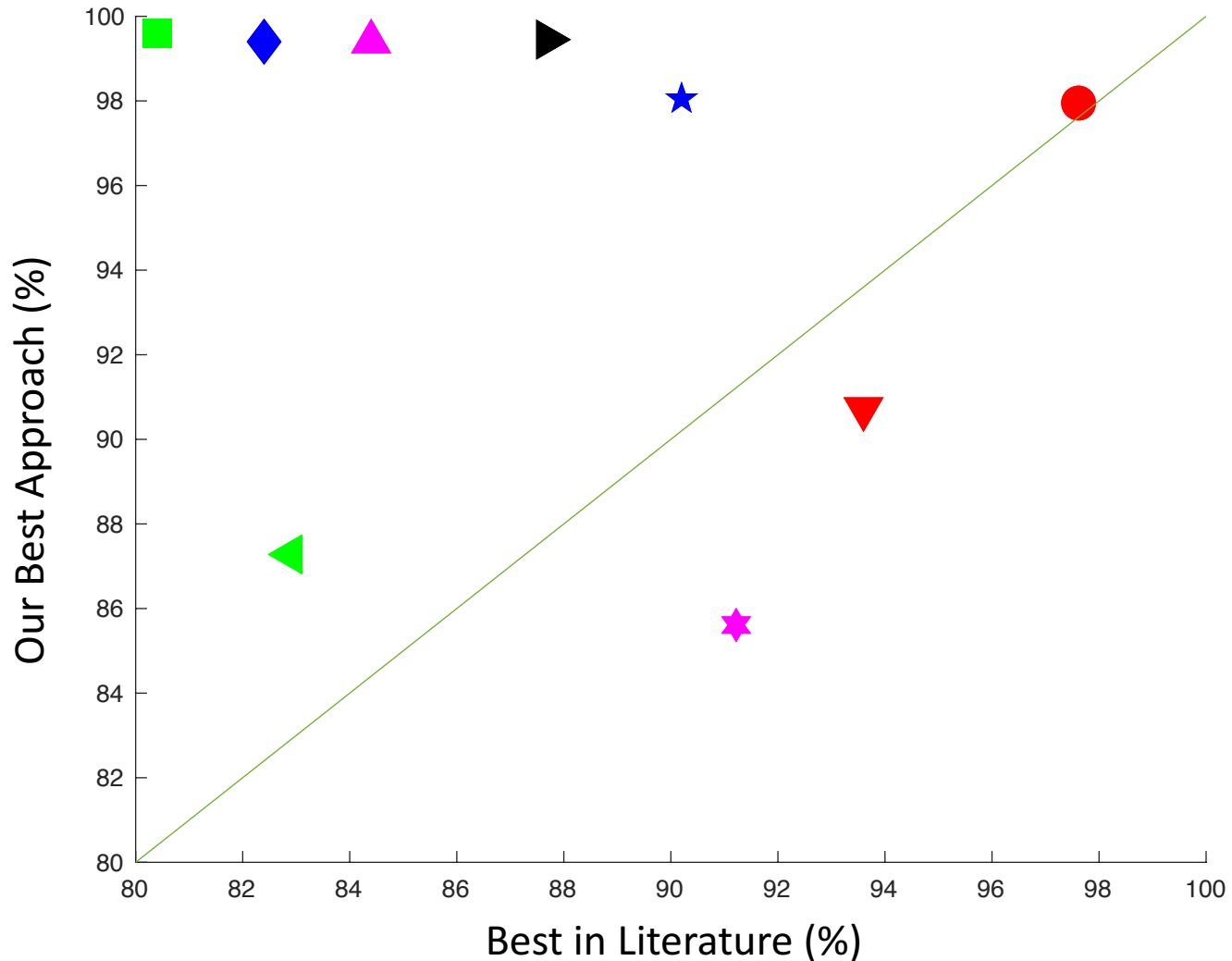


# Multiple Labels Example: Deep Learning for Sentiment Classification

- Train an anomaly detector on each polarity of sentiment
  - Can have many
  - E.g. positive & negative view
- Generate overall reconstruction error (ORE) for each label
- Use simple classifier on OREs to identify class



# Performance Analysis



Problem	Best Literature	Best Our Approach
● SMS Spam	97.64	<b>97.92</b>
■ MDSD-B(%)	80.4	<b>99.6</b>
◆ MDSD-D(%)	82.4	<b>99.4</b>
▲ MDSD-E(%)	84.4	<b>99.4</b>
▶ MDSD-K(%)	87.7	<b>99.45</b>
▼ Movie-Sub	<b>93.6</b>	90.72
◀ Movie-Rev1	82.9	<b>87.28</b>
★ Movie-Rev2	90.2	<b>98.05</b>
★ IMDB	<b>91.22</b>	85.61

# Summary



- Stacked Denoising Autoencoder allow us to spot anomalies within unlabeled data
- Probabilistic Topic Modelling allows us to look at documents at a higher level than just words
- When we have labels we can train on the labels to:
  - Identify sentiment
  - Classify the data

We Are recruiting:

- 2 PostDoc (Machine Learning / NLP)
- 1 PostDoc (Parallel Programming)
- Always looking for good PhD Candidates

[stephen.mcgough@newcastle.ac.uk](mailto:stephen.mcgough@newcastle.ac.uk)

[noura.al-moubayed@dur.ac.uk](mailto:noura.al-moubayed@dur.ac.uk)