# NLP Pre-Transformers

PD Dr. Juan J. Durillo

# Working with text

- Working with Neural Networks requires inputs in numerical representation
  - Character based representation (e.g., ascii code)
  - Word based encoding (each word a different representation)
    - Dictionary with all possible words; representation is based on the position on this dictionary

```
from tensorflow.keras.preprocessing.text import Tokenizer
sentences = ['Messi is the best player in the world', 'Barcelona
is the best team in the world']
tokenizer = Tokenizer(num_words=100)
tokenizer.fit_on_texts(sentences)
print(tokenizer.word_index)
{'the': 1, 'is': 2, 'best': 3, 'in': 4, 'world': 5, 'messi': 6,
'player': 7, 'barcelona': 8, 'team': 9}
```

# Working with text

- Alternative representation: One Hot Encoding
  - Vector of the dictionary length, with all components to 0 except 1
- Assuming the following dictionary

```
{'the': 1, 'is': 2, 'best': 3, 'in': 4, 'world': 5, 'messi':
6, 'player': 7, 'barcelona': 8, 'team': 9}
```
  - The word Messi would be represented by the vector

    `[0 0 0 0 0 1 0 0 0]`
  - The word player by

    `[0 0 0 0 0 0 1 0 0]`
  - The word the by

    `[1 0 0 0 0 0 0 0 0]`

# Text to Sequences

- A sequence (i.e., a sentence) is simply a list of (ordered) tokens
- Previous idea could be used for representing sentences

```
{'the': 1, 'is': 2, 'best': 3, 'in': 4, 'world': 5,
'messi': 6, 'player': 7, 'barcelona': 8, 'team': 9}
```

- The sentence 'Messi is the best player in the world' can be represented as the array

    [6, 2, 1, 3, 7, 4, 1, 5]

- And the sentence 'Barcelona is the best team in the world' can be represented as the array

    [8, 2, 1, 3, 9, 4, 1, 5]
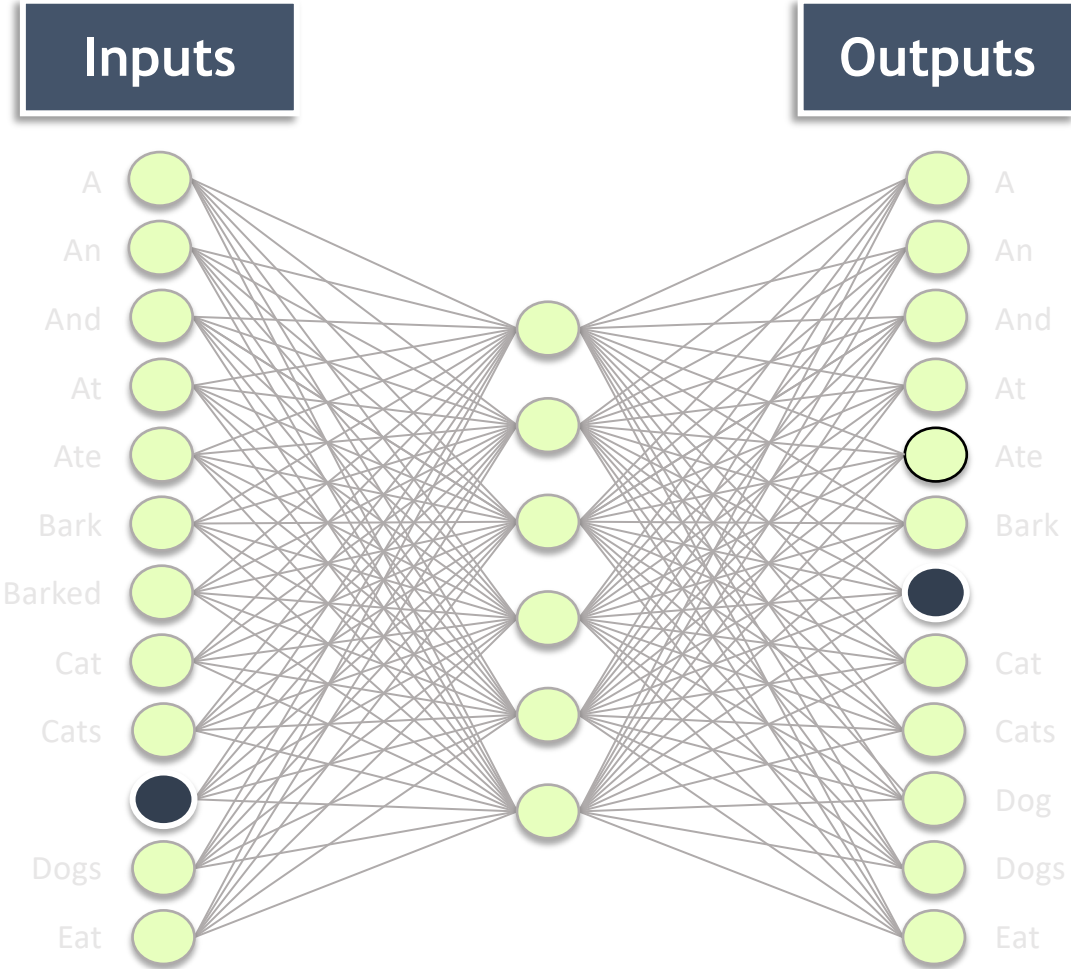
tokenizer.texts_to_sequences([str])

# Text to Sequences

- Alternatively, if the One Hot Encoding has been chosen, given the dictionary
  `{'the': 1, 'is': 2, 'best': 3, 'in': 4, 'world': 5, 'messi': 6, 'player': 7, 'barcelona': 8, 'team': 9}`
- The sentence 'Messi is the best player in the world' can be represented as the matrix

  ```
  [[0 0 0 0 0 1 0 0 0]
   [0 1 0 0 0 0 0 0 0]
   [1 0 0 0 0 0 0 0 0]
   [0 0 0 0 0 0 1 0 0]
   [0 0 0 1 0 0 0 0 0]
   [1 0 0 0 0 0 0 0 0]
   [0 0 0 0 1 0 0 0 0]]
  ```
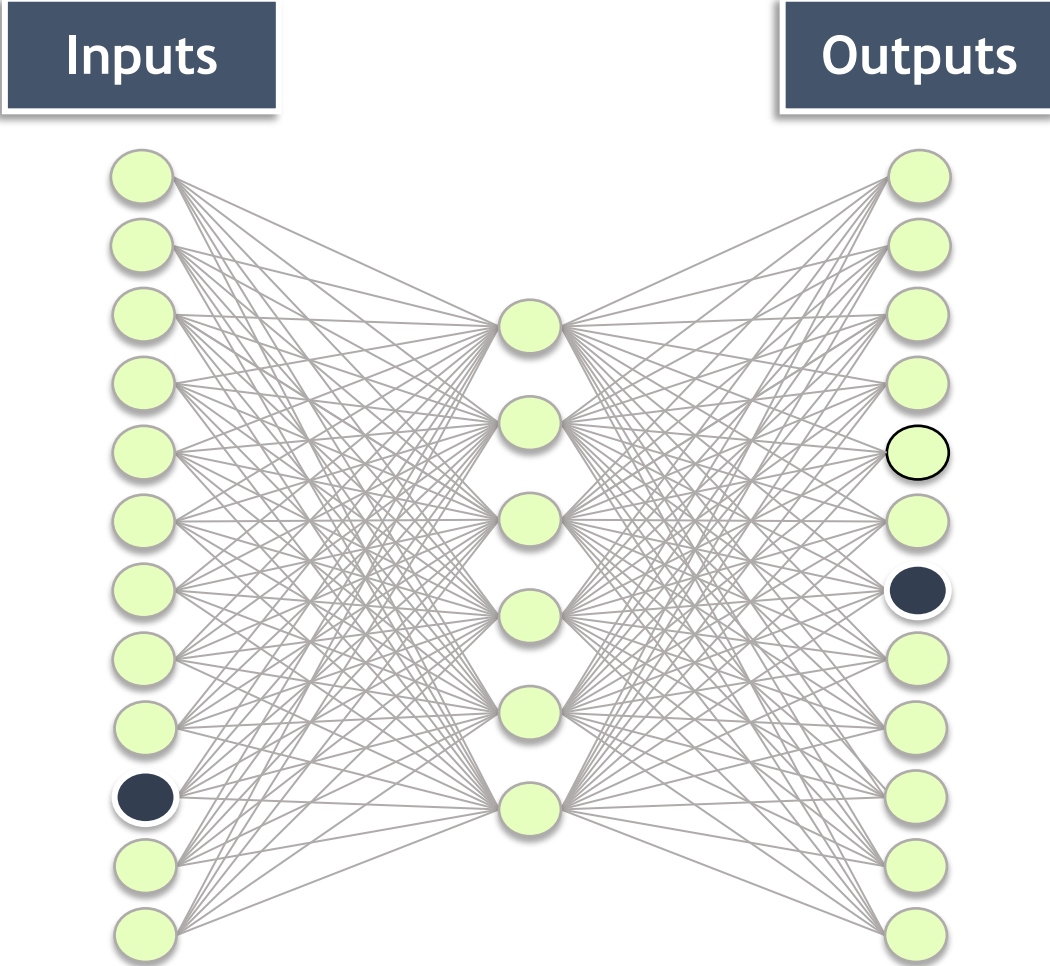
# From Words to Numbers

# From Words to Numbers



**Inputs**

**Outputs**

**Dictionary**

1. A
2. An
3. And
4. At
5. Ate
6. Bark
7. Barked

8. Cat
9. Cats
10. Dog
11. Dogs
12. Eat

# From Words to Numbers



Big

Giraffe
(.9, .9)

Llama
(-.9, .1)

Domestic                                                Wild

Falcon
(.15, -.4)

Kitty
(-.75, -.8)

Penguin
(.85, -.65)

Small

## Bigger Dictionary

| | | | | | | |
|---|---|---|---|---|---|---|
| 1. | A | 34. | Cat | 67. | An |
| 2. | An | 35. | Cats | 68. | And |
| 3. | And | 36. | Dog | 69. | At |
| 4. | At | 37. | Dogs | 70. | Ate |
| 5. | Ate | 38. | Eat | 71. | Bark |
| 6. | Bark | 39. | Eaten | 72. | Barked |
| 7. | Barked | 40. | A | 73. | Cat |
| 8. | Cat | 41. | An | 74. | Cats |
| 9. | Cats | 42. | And | 75. | Dog |
| 10. | Dog | 43. | At | 76. | Dogs |
| 11. | Dogs | 44. | Ate | 77. | Eat |
| 12. | Eat | 45. | Bark | 78. | Eaten |
| 13. | Eaten | 46. | Barked | 79. | … |
| 14. | A | 47. | Cat | 80. | … |
| 15. | An | 48. | Cats | 81. | … |
| 16. | And | 49. | Dog | 82. | … |
| 17. | At | 50. | Dogs | | |
| 18. | Ate | 51. | Eat | | |
| 19. | Bark | 52. | Eaten | | |
| 20. | Barked | 53. | A | | |
| 21. | Cat | 54. | An | | |
| 22. | Cats | 55. | And | | |
| 23. | Dog | 56. | At | | |
| 24. | Dogs | 57. | Ate | | |
| 25. | Eat | 58. | Bark | | |
| 26. | Eaten | 59. | Barked | | |
| 27. | A | 60. | Cat | | |
| 28. | An | 61. | Cats | | |
| 29. | And | 62. | Dog | | |
| 30. | At | 63. | Dogs | | |
| 31. | Ate | 64. | Eat | | |
| 32. | Bark | 65. | Eaten | | |
| 33. | Barked | 66. | A | | |

# From Words to Numbers



Inputs

Technically an Embedding

Outputs

Dictionary

1. A
2. An
3. And
4. At
5. Ate
6. Bark
7. Barked
8. Cat
9. Cats
10. Dog
11. Dogs
12. Eat

# Recurrent Neural Networks

# Learning From Text

- If you read the partial sentence:
  - Today there is an amazing blue …
- What do you think of next?

# Learning From Text

- If you read the partial sentence:
    - Today there is an amazing blue …

- What do you think of next?
    - Today there is an amazing blue sky.

# Learning from Text

- If you read the partial sentence:
    - She was born in Munich, therefore at school the primary language was ….
- In contrast to the previous example, the word that influences what we need to predict now is not the previous was, but was way beyond in the text
    - Do RNN still help in this case?

# Recurrent Neural Networks

"Cats say ____."

"Dogs say ____."

## Dictionary

1. Cats
2. Dogs
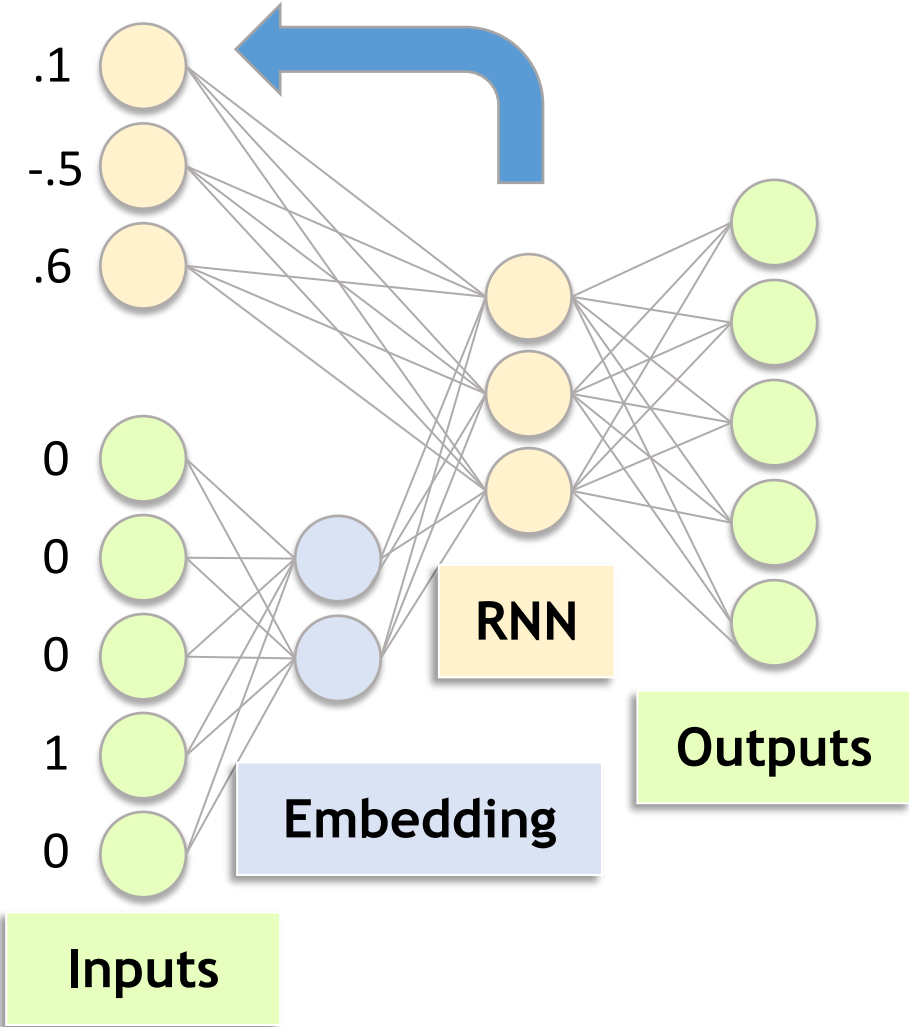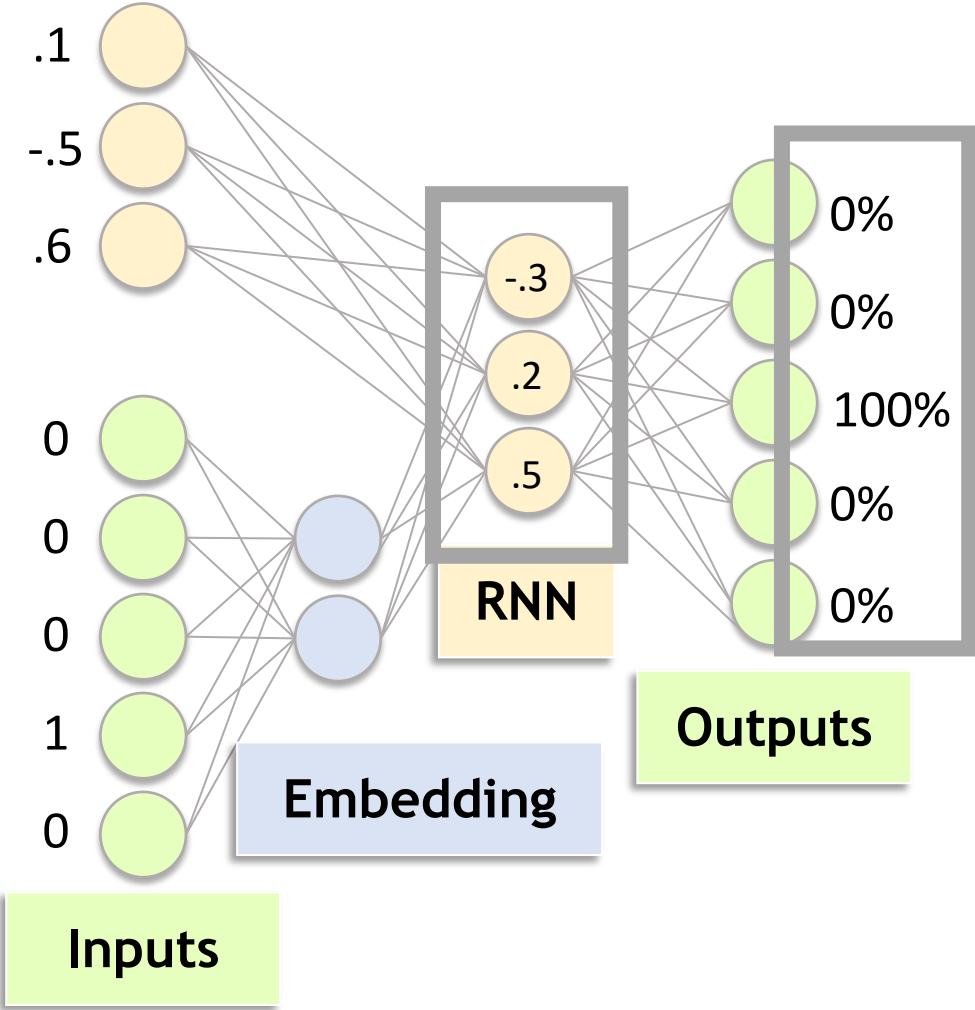3. Meow
4. Say
5. Woof

# Recurrent Neural Networks

# Recurrent Neural Networks

# Recurrent Neural Networks

# Recurrent Neural Networks



0

0

0

1

0

0

0

0

**Inputs**

**Embedding**

.1

-.5

.6

**RNN**

**Outputs**

"Cats say ____."

"Dogs say ____."

Dictionary

1. Cats
2. Dogs
3. Meow
4. Say
5. Woof

# Recurrent Neural Networks



.1
-.5
.6

0
0
0
1
0

**Inputs**

**Embedding**

**RNN**

**Outputs**

"Cats say ____."

"Dogs say ____."

Dictionary

1. Cats
2. Dogs
3. Meow
4. Say
5. Woof

# Recurrent Neural Networks

# Recurrent Neural Networks