# Building a Autoregressive Neural Network

**Part 1**

Luca WB

2026-01-28

# Table of contents

## 0.1 Brief summary

In this post, we will implement an Autoregressive Neural Network from scratch, relying solely on the PyTorch tensor class. We assume prior familiarity with Neural Networks; however, if your knowledge feels a bit rusty or you need a refresher, I recommend reading this post beforehand Building Neural Networks from Scratch.

The main reason for this is to learn how an Autoregressive NN works to generate words, for this, I'm drawing on Andrej Karpathy's video series about makemore, a network capable of creating more words of the same type, so if you train with names, it generates more proper names it generates more words that remember proper names, and so on with anything that is formed by letters.

In this post, I will cover how to make a simple model for our baseline, and how to implement a model with MLP and compare them.

## 0.2 Setup

First, you need to download PyTorch and the dataset. For PyTorch, just download in the official site https://pytorch.org/get-started/locally/. Now, for the dataset, you can create your own with random names that you can think, but It's much easier just download the names.txt dataset from the Andrej repository https://github.com/karpathy/makemore/blob/master/names.txt.

## 0.3 Creating a baseline

In propose of this, it's just to create the most simple and naive model. It's important because we need some baseline to compare with our future models, so we will create a model called bigram, the logic is just to look to the last character. Note that you will use just one character of context for our model, and we will consider that the most small part

of our word is a character, for models like chatGPT, they don't use characters, they use combinations of characters similar to syllables.

So, to start, we need first import our dataset and PyTorch

```
1  import torch
2
3  # Basically makes a list of all the names
4  names = open("names.txt", "r").read().splitlines()
5  names[:5]
```

```
['emma', 'olivia', 'ava', 'isabella', 'sophia']
```

Most part of models usually can't handle with characters, so it's useful to convert this letters in numbers in some way. For this, there are many possibles, but I will use just a simple dictionary to convert them. But we

```
1  chars = sorted(list(set("".join(names)))) # Creates an ordered list with all letters in ou
2  charToInt = {s:i+1 for i,s in enumerate(chars)} # Creates a dict to convert chars to int,
3  charToInt["."] = 0 # I will explain later why we need a special character
4  print(charToInt)
```

```
{'a': 1, 'b': 2, 'c': 3, 'd': 4, 'e': 5, 'f': 6, 'g': 7, 'h': 8, 'i': 9, 'j': 10, 'k': 11,
```

Just to get it ready, if we convert to int, so we can read it at the end, we will need an intToChar converter, so let's get it ready

```
1  intToChar = {s:i for i,s in charToInt.items()}
2  print(intToChar)
```

```
{1: 'a', 2: 'b', 3: 'c', 4: 'd', 5: 'e', 6: 'f', 7: 'g', 8: 'h', 9: 'i', 10: 'j', 11: 'k',
```

Know, for our model, we need to calculate the total number that each sequence occurs, like, with we start with letter "a", how many times occurs that "m" is the next character. And it's for this that we need and special characters, because we always need something to start, after all, the autoregressive model logic and take the output of the model and put it in its input, so we need an initial input. In our case, we will use "." as the symbol to start a name/words and to stop word (without a final symbol, it would generate forever). To make more clear, see the code bellow

```
1  N = torch.zeros((27,27)).int()
2
3  for name in names:
4      chars = ["."] + list(name) + ["."] # turn the name in a list of characters and add "."
5      for ch1,ch2 in zip(chars, chars[1:]): # In each loop, pick up one letter in ch1, and t
6          id1, id2 = charToInt[ch1], charToInt[ch2]
7          N[id1, id2] += 1
```

Basically, this count how often some sequence of characters occurs, like the most common letter sequence is "n" follow by ".", this mean, that the most commum letter to finish a name in our dataset it's "n". If you run with all the names, you can use the code bellow to find the most common occurrences

```
1  id1, id2 = (N == N.max()).nonzero(as_tuple=True) # Creates a boolean matrix that only it's
2  print(intToChar[id1.item()], "-->", intToChar[id2.item()],"occurs ", N.max().item())
```

```
n --> . occurs  6763
```

So let's see how our bigrams are distributed

```
1  import matplotlib.pyplot as plt
2
3  plt.figure(figsize= (16,16))
4  plt.imshow(N, cmap="Blues")
5  for i in range(27):
6      for j in range(27):
7          chstr = intToChar[i] + intToChar[j]
8          plt.text(j,i, chstr, ha="center", va="bottom", color="gray")
9          plt.text(j,i, N[i,j].item(), ha="center", va="top", color="gray")
10
11 plt.axis("off")
```

| | .. | .a | .b | .c | .d | .e | .f | .g | .h | .i | .j | .k | .l | .m | .n | .o | .p | .q | .r | .s | .t | .u | .v | .w | .x | .y | .z |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| . | 0 | 6440 | 1306 | 1542 | 1690 | 1531 | 417 | 669 | 874 | 591 | 2422 | 2963 | 1572 | 2538 | 1146 | 394 | 515 | 92 | 1639 | 2055 | 1308 | 78 | 376 | 307 | 134 | 535 | 929 |
| a | 6640 | 556 | 541 | 470 | 1042 | 692 | 134 | 168 | 2332 | 1650 | 175 | 568 | 2528 | 1634 | 5438 | 63 | 82 | 60 | 3264 | 1118 | 687 | 381 | 834 | 161 | 182 | 2050 | 435 |
| b | 114 | 321 | 38 | 1 | 65 | 655 | 0 | 0 | 41 | 217 | 1 | 0 | 103 | 0 | 4 | 105 | 0 | 0 | 842 | 8 | 2 | 45 | 0 | 0 | 0 | 83 | 0 |
| c | 97 | 815 | 0 | 42 | 1 | 551 | 0 | 2 | 664 | 271 | 3 | 316 | 116 | 0 | 0 | 380 | 1 | 11 | 76 | 5 | 35 | 35 | 0 | 3 | 0 | 104 | 4 |
| d | 516 | 1303 | 1 | 3 | 149 | 1283 | 5 | 25 | 118 | 674 | 9 | 3 | 60 | 30 | 31 | 378 | 0 | 1 | 424 | 29 | 4 | 92 | 17 | 23 | 0 | 317 | 1 |
| e | 3903 | 679 | 121 | 153 | 384 | 1271 | 82 | 125 | 152 | 818 | 55 | 178 | 3248 | 769 | 2675 | 83 | 83 | 14 | 1958 | 861 | 580 | 463 | 132 | 50 | 132 | 1070 | 181 |
| f | 80 | 242 | 0 | 0 | 0 | 123 | 44 | 1 | 1 | 160 | 0 | 2 | 20 | 0 | 4 | 60 | 0 | 0 | 114 | 6 | 18 | 10 | 0 | 4 | 0 | 14 | 2 |
| g | 108 | 330 | 3 | 0 | 19 | 334 | 1 | 25 | 360 | 190 | 3 | 0 | 32 | 6 | 27 | 83 | 0 | 0 | 201 | 30 | 31 | 85 | 1 | 26 | 0 | 31 | 1 |
| h | 2409 | 2244 | 8 | 2 | 24 | 674 | 2 | 2 | 1 | 729 | 9 | 29 | 185 | 117 | 138 | 287 | 1 | 1 | 204 | 31 | 71 | 166 | 39 | 10 | 0 | 213 | 20 |
| i | 2489 | 2445 | 110 | 509 | 440 | 1653 | 101 | 428 | 95 | 82 | 76 | 445 | 1345 | 427 | 2126 | 588 | 53 | 52 | 849 | 1316 | 541 | 269 | 269 | 8 | 89 | 779 | 277 |
| j | 71 | 1473 | 1 | 4 | 4 | 440 | 0 | 0 | 45 | 119 | 2 | 2 | 9 | 5 | 2 | 479 | 1 | 0 | 11 | 7 | 2 | 202 | 5 | 6 | 0 | 10 | 0 |
| k | 363 | 1731 | 2 | 2 | 2 | 895 | 1 | 0 | 307 | 509 | 2 | 0 | 139 | 9 | 26 | 344 | 0 | 0 | 109 | 95 | 17 | 50 | 2 | 34 | 0 | 379 | 2 |
| l | 1314 | 2623 | 52 | 25 | 138 | 2921 | 22 | 6 | 19 | 2480 | 6 | 24 | 1345 | 60 | 14 | 692 | 15 | 3 | 18 | 94 | 77 | 324 | 72 | 16 | 0 | 1588 | 10 |
| m | 516 | 2590 | 112 | 51 | 24 | 818 | 1 | 0 | 5 | 1256 | 7 | 1 | 5 | 168 | 20 | 452 | 38 | 0 | 97 | 35 | 4 | 139 | 3 | 2 | 0 | 287 | 11 |
| n | 6763 | 2977 | 8 | 213 | 704 | 1359 | 11 | 273 | 26 | 1725 | 44 | 58 | 195 | 19 | 1906 | 496 | 5 | 2 | 44 | 278 | 443 | 96 | 55 | 11 | 6 | 465 | 145 |
| o | 855 | 149 | 140 | 114 | 190 | 132 | 34 | 44 | 171 | 69 | 16 | 68 | 619 | 261 | 2411 | 115 | 95 | 3 | 1059 | 504 | 118 | 275 | 176 | 114 | 45 | 103 | 54 |
| p | 33 | 209 | 2 | 1 | 0 | 197 | 1 | 0 | 204 | 61 | 1 | 1 | 16 | 1 | 1 | 59 | 39 | 0 | 151 | 16 | 17 | 4 | 0 | 0 | 0 | 12 | 0 |
| q | 28 | 13 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 13 | 0 | 0 | 1 | 2 | 0 | 2 | 0 | 0 | 1 | 2 | 0 | 206 | 0 | 0 | 0 | 0 | 0 |
| r | 1377 | 2356 | 41 | 99 | 187 | 1697 | 9 | 76 | 121 | 3033 | 25 | 90 | 413 | 162 | 140 | 869 | 14 | 16 | 425 | 190 | 208 | 252 | 80 | 21 | 3 | 773 | 23 |
| s | 1169 | 1201 | 21 | 60 | 9 | 884 | 2 | 2 | 1285 | 684 | 2 | 82 | 279 | 90 | 24 | 531 | 51 | 1 | 55 | 461 | 765 | 185 | 14 | 24 | 0 | 215 | 10 |
| t | 483 | 1027 | 1 | 17 | 0 | 716 | 2 | 2 | 647 | 532 | 3 | 0 | 134 | 4 | 22 | 667 | 0 | 0 | 352 | 35 | 374 | 78 | 15 | 11 | 2 | 341 | 105 |
| u | 155 | 163 | 103 | 103 | 136 | 169 | 19 | 47 | 58 | 121 | 14 | 93 | 301 | 154 | 275 | 16 | 16 | 10 | 414 | 474 | 82 | 3 | 37 | 86 | 34 | 13 | 45 |
| v | 88 | 642 | 1 | 0 | 1 | 568 | 0 | 0 | 1 | 911 | 0 | 3 | 14 | 0 | 8 | 153 | 0 | 0 | 48 | 0 | 0 | 7 | 7 | 0 | 0 | 121 | 0 |
| w | 51 | 280 | 1 | 0 | 8 | 149 | 2 | 1 | 23 | 192 | 0 | 0 | 13 | 0 | 58 | 36 | 0 | 0 | 22 | 20 | 8 | 25 | 0 | 2 | 0 | 73 | 1 |
| x | 164 | 103 | 1 | 4 | 5 | 36 | 3 | 0 | 1 | 102 | 0 | 0 | 39 | 1 | 1 | 41 | 0 | 0 | 0 | 31 | 70 | 5 | 0 | 3 | 38 | 30 | 19 |
| y | 2007 | 2143 | 27 | 115 | 272 | 301 | 12 | 30 | 22 | 192 | 23 | 86 | 1104 | 148 | 1826 | 271 | 15 | 6 | 291 | 401 | 104 | 141 | 106 | 4 | 28 | 23 | 78 |
| z | 160 | 860 | 4 | 2 | 2 | 373 | 0 | 1 | 43 | 364 | 2 | 2 | 123 | 35 | 4 | 110 | 2 | 0 | 32 | 4 | 4 | 73 | 2 | 3 | 1 | 147 | 45 |

One thing very interesting you can note, it's that have many combinations that don't exist, like "bk" or "gc". This makes it impossible for our model to generate a name with this combination, it is ok to leave it like this, but it would be a good practice to add 1 in all values, thus ensuring that at least there is the minimal possibility of generating a rare sequence

```
1   N = N + 1
```

| | . | a | b | c | d | e | f | g | h | i | j | k | l | m | n | o | p | q | r | s | t | u | v | w | x | y | z |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| . | 1 | | 1307 | 1543 | 1691 | 1532 | 418 | 670 | 875 | 592 | 2423 | 2964 | 1573 | 2539 | 1147 | 395 | 516 | 93 | 1640 | 2056 | 1309 | 79 | 377 | 308 | 135 | 536 | 930 |
| a | 6641 | 557 | 542 | 471 | 1043 | 693 | 135 | 169 | 2333 | 1651 | 176 | 569 | 2529 | 1635 | 5439 | 64 | 83 | 61 | 3265 | 1119 | 688 | 382 | 835 | 162 | 183 | 2051 | 436 |
| b | 115 | 322 | 39 | 2 | 66 | 656 | 1 | 1 | 42 | 218 | 2 | 1 | 104 | 1 | 5 | 106 | 1 | 1 | 843 | 9 | 3 | 46 | 1 | 1 | 1 | 84 | 1 |
| c | 98 | 816 | 1 | 43 | 2 | 552 | 1 | 3 | 665 | 272 | 4 | 317 | 117 | 1 | 1 | 381 | 2 | 12 | 77 | 6 | 36 | 36 | 1 | 4 | 1 | 105 | 5 |
| d | 517 | 1304 | 2 | 4 | 150 | 1284 | 6 | 26 | 119 | 675 | 10 | 4 | 61 | 31 | 32 | 379 | 1 | 2 | 425 | 30 | 5 | 93 | 18 | 24 | 1 | 318 | 2 |
| e | 3984 | 680 | 122 | 154 | 385 | 1272 | 83 | 126 | 153 | 819 | 56 | 179 | 3249 | 770 | 2676 | 84 | 15 | | 1959 | 862 | 581 | 70 | 464 | 133 | 133 | 1071 | 182 |
| f | 81 | 243 | 1 | 1 | 1 | 124 | 45 | 2 | 2 | 161 | 1 | 3 | 21 | 1 | 5 | 61 | 1 | 1 | 115 | 7 | 19 | 11 | 1 | 5 | 1 | 15 | 3 |
| g | 109 | 331 | 4 | 1 | 20 | 335 | 2 | 26 | 361 | 191 | 4 | 1 | 33 | 7 | 28 | 84 | 1 | 1 | 202 | 31 | 32 | 86 | 2 | 27 | 1 | 32 | 2 |
| h | 2410 | 2245 | 9 | 3 | 25 | 675 | 3 | 3 | 2 | 730 | 10 | 30 | 186 | 118 | 139 | 288 | 2 | 2 | 205 | 32 | 72 | 167 | 40 | 11 | 1 | 214 | 21 |
| i | 2490 | 2446 | 111 | 510 | 441 | 1654 | 102 | 429 | 96 | 83 | 77 | 446 | 1346 | 428 | 2127 | 589 | 54 | 53 | 850 | 1317 | 542 | 110 | 270 | 9 | 90 | 780 | 278 |
| j | 72 | 1474 | 2 | 5 | 5 | 441 | 1 | 1 | 46 | 120 | 3 | 3 | 10 | 6 | 3 | 480 | 2 | 1 | 12 | 8 | 3 | 203 | 6 | 7 | 1 | 11 | 1 |
| k | 364 | 1732 | 3 | 3 | 3 | 896 | 2 | 1 | 308 | 510 | 3 | 21 | 140 | 10 | 27 | 345 | 1 | 1 | 110 | 96 | 18 | 51 | 3 | 35 | 1 | 380 | 3 |
| l | 1315 | 2624 | 53 | 26 | 139 | 2922 | 23 | 7 | 20 | 2481 | 7 | 25 | 1346 | 61 | 15 | 693 | 16 | 4 | 19 | 95 | 78 | 325 | 73 | 17 | 1 | 1589 | 11 |
| m | 517 | 2591 | 113 | 52 | 25 | 819 | 2 | 1 | 6 | 1257 | 8 | 2 | 6 | 169 | 21 | 453 | 39 | 1 | 98 | 36 | 5 | 140 | 4 | 3 | 1 | 288 | 12 |
| n | 6764 | 2978 | 9 | 214 | 705 | 1360 | 12 | 274 | 27 | 1726 | 45 | 59 | 196 | 20 | 1907 | 497 | 6 | 3 | 45 | 279 | 444 | 97 | 56 | 12 | 7 | 466 | 146 |
| o | 856 | 150 | 141 | 115 | 191 | 133 | 35 | 45 | 172 | 70 | 17 | 69 | 620 | 262 | 2412 | 116 | 96 | 4 | 1060 | 505 | 119 | 276 | 177 | 115 | 46 | 104 | 55 |
| p | 34 | 210 | 3 | 2 | 1 | 198 | 2 | 1 | 205 | 62 | 2 | 2 | 17 | 2 | 2 | 60 | 40 | 2 | 152 | 17 | 18 | 5 | 1 | 1 | 1 | 13 | 1 |
| q | 29 | 14 | 1 | 1 | 1 | 2 | 1 | 1 | 1 | 14 | 1 | 1 | 2 | 3 | 1 | 3 | 1 | 1 | 2 | 1 | 1 | 207 | 1 | 1 | 1 | 1 | 1 |
| r | 1378 | 2357 | 42 | 100 | 188 | 1698 | 10 | 77 | 122 | 3034 | 26 | 91 | 414 | 163 | 141 | 870 | 15 | 17 | 426 | 191 | 209 | 253 | 81 | 22 | 4 | 774 | 24 |
| s | 1170 | 1202 | 22 | 61 | 10 | 885 | 3 | 3 | 1286 | 685 | 3 | 83 | 280 | 91 | 25 | 532 | 52 | 2 | 56 | 462 | 766 | 186 | 15 | 25 | 1 | 216 | 11 |
| t | 484 | 1028 | 2 | 18 | 1 | 717 | 3 | 3 | 648 | 533 | 4 | 1 | 135 | 5 | 23 | 668 | 1 | 1 | 353 | 36 | 375 | 79 | 16 | 12 | 3 | 342 | 106 |
| u | 156 | 164 | 104 | 104 | 137 | 170 | 20 | 48 | 59 | 122 | 15 | 94 | 302 | 155 | 276 | 10 | 17 | 11 | 415 | 475 | 83 | 4 | 38 | 87 | 35 | 14 | 46 |
| v | 89 | 643 | 2 | 1 | 2 | 569 | 1 | 1 | 2 | 912 | 1 | 4 | 15 | 1 | 9 | 154 | 1 | 1 | 49 | 1 | 1 | 8 | 8 | 1 | 1 | 122 | 1 |
| w | 52 | 281 | 2 | 1 | 9 | 150 | 3 | 2 | 24 | 149 | 1 | 7 | 14 | 3 | 59 | 37 | 1 | 1 | 23 | 21 | 9 | 26 | 1 | 3 | 1 | 74 | 2 |
| x | 165 | 104 | 2 | 5 | 6 | 37 | 4 | 1 | 2 | 103 | 1 | 1 | 40 | 2 | 2 | 42 | 1 | 1 | 1 | 32 | 71 | 6 | 1 | 4 | 39 | 31 | 20 |
| y | 2008 | 2144 | 28 | 116 | 273 | 302 | 13 | 31 | 23 | 193 | 24 | 87 | 1105 | 149 | 1827 | 272 | 16 | 7 | 292 | 402 | 105 | 142 | 107 | 5 | 29 | 24 | 79 |
| z | 161 | 861 | 5 | 3 | 3 | 374 | 1 | 44 | 2 | 365 | 3 | 3 | 124 | 36 | 5 | 111 | 3 | 1 | 33 | 5 | 5 | 74 | 3 | 4 | 2 | 148 | 46 |

So, lets transform our probability matrix

```python
P = N
P = P / P.sum(dim=1, keepdims=True)
```
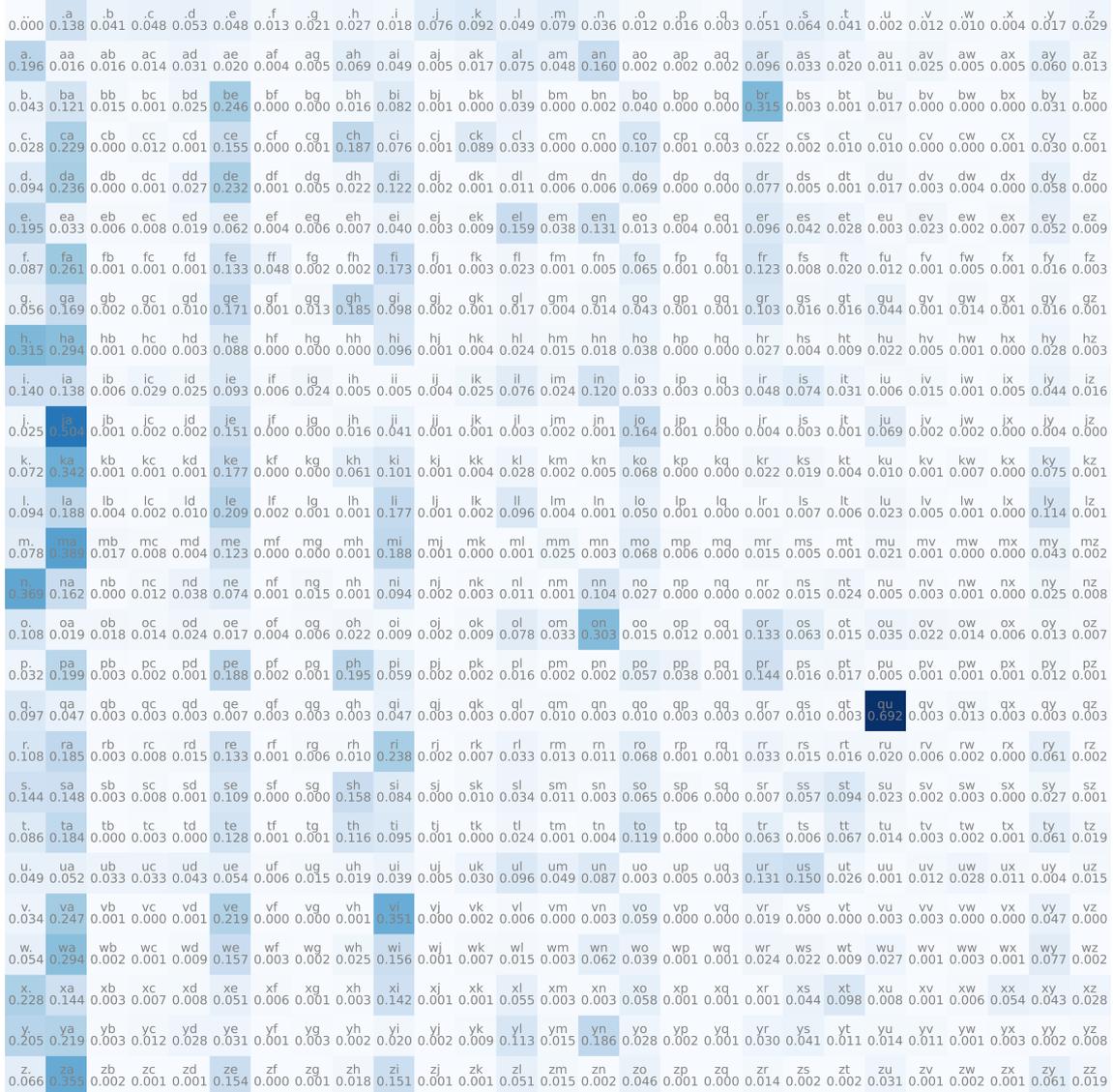
```python
import matplotlib.pyplot as plt

plt.figure(figsize= (16,16))
plt.imshow(P, cmap="Blues")
for i in range(27):
    for j in range(27):
        chstr = intToChar[i] + intToChar[j]
        plt.text(j,i, chstr, ha="center", va="bottom", color="gray")
```

```python
        plt.text(j,i, f"{P[i,j].item():.3f}", ha="center", va="top", color="gray")

plt.axis("off")
```

| | . | a | b | c | d | e | f | g | h | i | j | k | l | m | n | o | p | q | r | s | t | u | v | w | x | y | z |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| . | 0.000 | 0.138 | 0.041 | 0.048 | 0.053 | 0.048 | 0.013 | 0.021 | 0.027 | 0.018 | 0.076 | 0.092 | 0.049 | 0.079 | 0.036 | 0.012 | 0.016 | 0.003 | 0.051 | 0.064 | 0.041 | 0.002 | 0.012 | 0.010 | 0.004 | 0.017 | 0.029 |
| a | 0.196 | 0.016 | 0.016 | 0.014 | 0.031 | 0.020 | 0.004 | 0.005 | 0.069 | 0.049 | 0.005 | 0.017 | 0.075 | 0.048 | 0.160 | 0.002 | 0.002 | 0.002 | 0.096 | 0.033 | 0.020 | 0.011 | 0.025 | 0.005 | 0.005 | 0.060 | 0.013 |
| b | 0.043 | 0.121 | 0.015 | 0.001 | 0.025 | 0.246 | 0.000 | 0.000 | 0.016 | 0.082 | 0.001 | 0.000 | 0.039 | 0.000 | 0.002 | 0.040 | 0.000 | 0.000 | 0.315 | 0.003 | 0.001 | 0.017 | 0.000 | 0.000 | 0.000 | 0.031 | 0.000 |
| c | 0.028 | 0.229 | 0.000 | 0.012 | 0.001 | 0.155 | 0.000 | 0.001 | 0.187 | 0.076 | 0.001 | 0.089 | 0.033 | 0.000 | 0.000 | 0.107 | 0.001 | 0.003 | 0.022 | 0.002 | 0.010 | 0.010 | 0.000 | 0.000 | 0.001 | 0.030 | 0.001 |
| d | 0.094 | 0.236 | 0.000 | 0.001 | 0.027 | 0.232 | 0.001 | 0.005 | 0.022 | 0.122 | 0.002 | 0.001 | 0.011 | 0.006 | 0.006 | 0.069 | 0.000 | 0.000 | 0.077 | 0.005 | 0.001 | 0.017 | 0.003 | 0.004 | 0.000 | 0.058 | 0.000 |
| e | 0.195 | 0.033 | 0.006 | 0.008 | 0.019 | 0.062 | 0.004 | 0.006 | 0.007 | 0.040 | 0.003 | 0.009 | 0.159 | 0.038 | 0.131 | 0.013 | 0.004 | 0.001 | 0.096 | 0.042 | 0.028 | 0.003 | 0.023 | 0.002 | 0.007 | 0.052 | 0.009 |
| f | 0.087 | 0.261 | 0.001 | 0.001 | 0.001 | 0.133 | 0.048 | 0.002 | 0.002 | 0.173 | 0.001 | 0.003 | 0.023 | 0.001 | 0.005 | 0.065 | 0.001 | 0.001 | 0.123 | 0.008 | 0.020 | 0.012 | 0.001 | 0.005 | 0.001 | 0.016 | 0.003 |
| g | 0.056 | 0.169 | 0.002 | 0.001 | 0.010 | 0.171 | 0.001 | 0.013 | 0.185 | 0.098 | 0.002 | 0.001 | 0.017 | 0.004 | 0.014 | 0.043 | 0.001 | 0.001 | 0.103 | 0.016 | 0.016 | 0.044 | 0.001 | 0.014 | 0.001 | 0.016 | 0.001 |
| h | 0.315 | 0.294 | 0.001 | 0.000 | 0.003 | 0.088 | 0.000 | 0.000 | 0.000 | 0.096 | 0.001 | 0.004 | 0.024 | 0.015 | 0.018 | 0.038 | 0.000 | 0.000 | 0.027 | 0.004 | 0.009 | 0.022 | 0.005 | 0.001 | 0.000 | 0.028 | 0.003 |
| i | 0.140 | 0.138 | 0.006 | 0.029 | 0.025 | 0.093 | 0.006 | 0.024 | 0.005 | 0.005 | 0.004 | 0.025 | 0.076 | 0.024 | 0.120 | 0.033 | 0.003 | 0.003 | 0.048 | 0.074 | 0.031 | 0.006 | 0.015 | 0.001 | 0.005 | 0.044 | 0.016 |
| j | 0.025 | 0.504 | 0.001 | 0.002 | 0.002 | 0.151 | 0.000 | 0.000 | 0.016 | 0.041 | 0.001 | 0.001 | 0.003 | 0.002 | 0.001 | 0.164 | 0.001 | 0.000 | 0.004 | 0.003 | 0.001 | 0.069 | 0.002 | 0.002 | 0.000 | 0.004 | 0.000 |
| k | 0.072 | 0.342 | 0.001 | 0.001 | 0.001 | 0.177 | 0.000 | 0.000 | 0.061 | 0.101 | 0.001 | 0.004 | 0.028 | 0.002 | 0.005 | 0.068 | 0.000 | 0.000 | 0.022 | 0.019 | 0.004 | 0.010 | 0.001 | 0.007 | 0.000 | 0.075 | 0.001 |
| l | 0.094 | 0.188 | 0.004 | 0.002 | 0.010 | 0.209 | 0.002 | 0.001 | 0.001 | 0.177 | 0.001 | 0.002 | 0.096 | 0.004 | 0.001 | 0.050 | 0.001 | 0.000 | 0.001 | 0.007 | 0.006 | 0.023 | 0.005 | 0.001 | 0.001 | 0.114 | 0.001 |
| m | 0.078 | 0.389 | 0.017 | 0.008 | 0.004 | 0.123 | 0.000 | 0.000 | 0.001 | 0.188 | 0.001 | 0.000 | 0.001 | 0.025 | 0.003 | 0.068 | 0.006 | 0.000 | 0.015 | 0.005 | 0.001 | 0.021 | 0.001 | 0.000 | 0.000 | 0.043 | 0.002 |
| n | 0.369 | 0.162 | 0.000 | 0.012 | 0.038 | 0.074 | 0.001 | 0.015 | 0.001 | 0.094 | 0.002 | 0.003 | 0.011 | 0.001 | 0.104 | 0.027 | 0.000 | 0.000 | 0.002 | 0.015 | 0.024 | 0.005 | 0.003 | 0.001 | 0.000 | 0.025 | 0.008 |
| o | 0.108 | 0.019 | 0.018 | 0.014 | 0.024 | 0.017 | 0.004 | 0.006 | 0.022 | 0.009 | 0.002 | 0.009 | 0.078 | 0.033 | 0.303 | 0.015 | 0.012 | 0.001 | 0.133 | 0.063 | 0.015 | 0.035 | 0.022 | 0.014 | 0.006 | 0.013 | 0.007 |
| p | 0.032 | 0.199 | 0.003 | 0.002 | 0.001 | 0.188 | 0.002 | 0.001 | 0.195 | 0.059 | 0.002 | 0.002 | 0.016 | 0.002 | 0.002 | 0.057 | 0.038 | 0.001 | 0.144 | 0.016 | 0.017 | 0.005 | 0.001 | 0.001 | 0.001 | 0.012 | 0.001 |
| q | 0.097 | 0.047 | 0.003 | 0.003 | 0.003 | 0.007 | 0.003 | 0.003 | 0.003 | 0.047 | 0.003 | 0.003 | 0.003 | 0.007 | 0.010 | 0.003 | 0.010 | 0.003 | 0.007 | 0.010 | 0.003 | 0.692 | 0.003 | 0.013 | 0.003 | 0.003 | 0.003 |
| r | 0.108 | 0.185 | 0.003 | 0.008 | 0.015 | 0.133 | 0.001 | 0.006 | 0.010 | 0.238 | 0.002 | 0.007 | 0.033 | 0.013 | 0.011 | 0.068 | 0.001 | 0.001 | 0.033 | 0.015 | 0.016 | 0.020 | 0.006 | 0.002 | 0.000 | 0.061 | 0.002 |
| s | 0.144 | 0.148 | 0.003 | 0.008 | 0.001 | 0.109 | 0.000 | 0.000 | 0.158 | 0.084 | 0.000 | 0.010 | 0.034 | 0.011 | 0.003 | 0.065 | 0.006 | 0.000 | 0.007 | 0.057 | 0.094 | 0.023 | 0.002 | 0.003 | 0.000 | 0.027 | 0.001 |
| t | 0.086 | 0.184 | 0.000 | 0.003 | 0.000 | 0.128 | 0.001 | 0.001 | 0.116 | 0.095 | 0.001 | 0.000 | 0.024 | 0.001 | 0.004 | 0.119 | 0.000 | 0.000 | 0.063 | 0.006 | 0.067 | 0.014 | 0.003 | 0.002 | 0.001 | 0.061 | 0.019 |
| u | 0.049 | 0.052 | 0.033 | 0.033 | 0.043 | 0.054 | 0.006 | 0.015 | 0.019 | 0.039 | 0.005 | 0.030 | 0.096 | 0.049 | 0.087 | 0.003 | 0.005 | 0.003 | 0.131 | 0.150 | 0.026 | 0.001 | 0.012 | 0.028 | 0.011 | 0.004 | 0.015 |
| v | 0.034 | 0.247 | 0.001 | 0.000 | 0.001 | 0.219 | 0.000 | 0.000 | 0.001 | 0.351 | 0.000 | 0.002 | 0.006 | 0.000 | 0.003 | 0.059 | 0.000 | 0.000 | 0.019 | 0.000 | 0.000 | 0.003 | 0.003 | 0.000 | 0.000 | 0.047 | 0.000 |
| w | 0.054 | 0.294 | 0.002 | 0.001 | 0.009 | 0.157 | 0.003 | 0.002 | 0.025 | 0.156 | 0.001 | 0.007 | 0.015 | 0.003 | 0.062 | 0.039 | 0.001 | 0.001 | 0.024 | 0.022 | 0.009 | 0.027 | 0.001 | 0.003 | 0.001 | 0.077 | 0.002 |
| x | 0.228 | 0.144 | 0.003 | 0.007 | 0.008 | 0.051 | 0.006 | 0.001 | 0.003 | 0.142 | 0.001 | 0.001 | 0.055 | 0.003 | 0.003 | 0.058 | 0.001 | 0.001 | 0.001 | 0.044 | 0.098 | 0.008 | 0.001 | 0.006 | 0.054 | 0.043 | 0.028 |
| y | 0.205 | 0.219 | 0.003 | 0.012 | 0.028 | 0.031 | 0.001 | 0.003 | 0.002 | 0.020 | 0.002 | 0.009 | 0.113 | 0.015 | 0.186 | 0.028 | 0.002 | 0.001 | 0.030 | 0.041 | 0.011 | 0.014 | 0.011 | 0.001 | 0.003 | 0.002 | 0.008 |
| z | 0.066 | 0.355 | 0.002 | 0.001 | 0.001 | 0.154 | 0.000 | 0.001 | 0.018 | 0.151 | 0.001 | 0.001 | 0.051 | 0.015 | 0.002 | 0.046 | 0.001 | 0.000 | 0.014 | 0.002 | 0.002 | 0.031 | 0.001 | 0.002 | 0.001 | 0.061 | 0.019 |

Some probabilities stay in 0 because the visualization it's limited to 3 decimal numbers. Now we already have our model, it's just our probability matrix P, bellow I will show how to use it.

```python
import random

for i in range(10):
    out = []
```

```
5      init = 0
6      while True:
7          id = torch.multinomial(P[init], num_samples=1, replacement=True).item()
8
9          if id == 0:
10             break
11
12         out.append(intToChar[id])
13         init = id
14     print("".join(out))
```

```
ladhai
ken
ile
chiliariali
jamitt
janany
h
sekal
trlelerilynemi
llsdrgaabr
```