# 📘 N-gram Language Model — Full Theory + Implementation (with Interpolation)

```python
1  # Import the random module for generating random numbers,
2  # which can be useful for sampling or shuffling tokens.
3  import random
4
5  # Import the math module for mathematical operations
6  # (e.g., logarithms, exponentiation, probabilities).
7  import math
8
9  # Import Counter from collections to easily count the frequency
10 # of n-grams and contexts in the language model.
11 from collections import Counter
12
13 # Import List, Tuple, and Dict from typing to provide
14 # type hints for function parameters and return values.
15 from typing import List, Tuple, Dict
16
```

## ◆ Theory

Before we can model language statistically, we must represent text as discrete **tokens** (usually words). Tokenization is the process of splitting sentences into individual words or symbols. Adding special tokens:

- `<s>` marks the **start of a sentence**
- `</s>` marks the **end of a sentence**

This helps the model learn context boundaries, so it doesn't mix tokens from different sentences.

## ◆ Function

## 🧩 1. Tokenization and Preprocessing

```python
1  def tokenize_text(text: str) -> List[str]:
2      """
3      Theory:
4      Converts raw text into tokens that the language model can process.
5      It splits text into sentences, lowercases words for normalization,
6      and inserts start (<s>) and end (</s>) markers around each sentence.
7
8      Practical importance:
9      - Enables the model to learn where sentences begin and end.
10     - Normalizes text to reduce vocabulary sparsity.
11     """
12     # Split the input text into sentences using '.' as a delimiter
13     # and remove any leading/trailing spaces from each sentence
14     sentences = [s.strip() for s in text.split('.') if s.strip()]
15
16     # Initialize an empty list to hold all tokens
17     tokens = []
18
19     # Loop through each sentence in the text
20     for s in sentences:
21         # Convert the sentence to lowercase for normalization
22         words = s.lower().split()
23
24         # Add special start and end markers around each sentence
25         # <s> indicates the start of a sentence, </s> indicates the end
26         tokens.extend(["<s>"] + words + ["</s>"])
27
28     # Return the complete list of tokens
29     return tokens
```

```python
1  tokens = tokenize_text("This is a sample sentence. It has multiple words.")
2  tokens
```

```
['<s>',
 'this',
 'is',
 'a',
 'sample',
 'sentence',
 '</s>',
 '<s>',
 'it',
 'has',
 'multiple',
 'words',
 '</s>']
```

## ✳️ 2. Building N-grams

### ◆ Theory

An **N-gram** is a contiguous sequence of (N) items (usually words). For example, for "I love AI":

- Unigrams (N=1): I, love, AI
- Bigrams (N=2): (I, love), (love, AI)
- Trigrams (N=3): (I, love, AI)

These sequences capture **local word dependencies** — how likely a word is to follow a specific context.

### ◆ Function

```python
1  def build_ngrams(tokens: List[str], n: int) -> List[Tuple[str, ...]]:
2      """
3      Theory:
4      Generates n-grams from a list of tokens.
5      Each n-gram represents a potential local context used for probability estimation.
6
7      Example:
8          tokens = ["I", "love", "AI"]
9          n=2 → [("I","love"), ("love","AI")]
10
11      Purpose:
12      - Helps us measure how often specific word patterns occur in the corpus.
13      """
14      return [tuple(tokens[i:i+n]) for i in range(len(tokens) - n + 1)]
```

```python
1  build_ngrams(tokens, 2)
```

```
[('<s>', 'this'),
 ('this', 'is'),
 ('is', 'a'),
 ('a', 'sample'),
 ('sample', 'sentence'),
 ('sentence', '</s>'),
 ('</s>', '<s>'),
 ('<s>', 'it'),
 ('it', 'has'),
 ('has', 'multiple'),
 ('multiple', 'words'),
 ('words', '</s>')]
```

```python
1  build_ngrams(tokens, 3)
```

```
[('<s>', 'this', 'is'),
 ('this', 'is', 'a'),
 ('is', 'a', 'sample'),
 ('a', 'sample', 'sentence'),
 ('sample', 'sentence', '</s>'),
 ('sentence', '</s>', '<s>'),
 ('</s>', '<s>', 'it'),
 ('<s>', 'it', 'has'),
 ('it', 'has', 'multiple'),
 ('has', 'multiple', 'words'),
 ('multiple', 'words', '</s>')]
```

## 🧩 3. Training the Model (Counting Frequencies)

### 🔷 Theory

Language modeling relies on estimating **conditional probabilities**:

$$P(w_i|w_{i-n+1}, ..., w_{i-1})$$

This is estimated via **Maximum Likelihood Estimation (MLE)**:

$$P(w_i|context) = \frac{Count(context, w_i)}{Count(context)}$$

Thus, we need to count:

- How often each **n-gram** appears.
- How often each **context** (n-1-gram) appears.

### 🔷 Function

```python
1 def train_ngram_model(tokens: List[str], n: int) -> Tuple[Dict[Tuple[str, ...], int], Dict[Tuple[str, ...], int]]:
2     """
3     Theory:
4     - Collect frequency counts for n-grams and their (n-1)-gram contexts.
5     - These counts are the empirical statistics from which we estimate probabilities.
6
7     Importance:
8     - The counts form the 'memory' of our language model.
9     - They are later used for both smoothing and interpolation.
10
11    Returns:
12        ngram_counts   → frequency of each n-gram
13        context_counts → frequency of each (n-1)-gram context
14    """
15    # Step 1: Build all n-grams from the tokenized text
16    # Example for n=3: ("<s>", "i", "am"), ("i", "am", "happy"), ...
17    # Count how many times each unique n-gram appears
18    ngram_counts = Counter(build_ngrams(tokens, n))
19
20    # Step 2: Build all (n-1)-grams (contexts for n-grams)
21    # Example for n=3, context = ("<s>", "i"), ("i", "am"), ...
22    # These context counts help compute conditional probabilities later
23    context_counts = Counter(build_ngrams(tokens, n - 1))
24
25    # Step 3: Return both dictionaries:
26    # - ngram_counts: frequency of each n-gram
27    # - context_counts: frequency of each (n-1)-gram
28    return ngram_counts, context_counts
```

```python
1 ngram_counts, context_counts = train_ngram_model(tokens, 2)
```

## 🧩 4. Add-k (Laplace) Smoothing

### ◆ Theory

Raw MLE assigns **zero probability** to unseen n-grams, which breaks models — especially during generation or perplexity evaluation.

**Add-k smoothing** fixes this by pretending we've seen every n-gram a small number of times $k$.

Formula:

$$P_{smooth}(w_i|context) = \frac{Count(context, w_i) + k}{Count(context) + k|V|}$$

where $|V|$ is vocabulary size.

This ensures **no zero probabilities** and maintains normalization.

### ◆ Function

```python
def smoothed_probability(
    ngram: Tuple[str, ...],
    ngram_counts: Dict,
    context_counts: Dict,
    vocab_size: int,
    k: float = 1.0
) -> float:
    """
    Computes the smoothed probability of an n-gram using Laplace (Add-k) smoothing.
    """

    # Extract the context (all words except the last one)
    # Example: if ngram = ("i", "love"), then context = ("i",)
    context = ngram[:-1]

    # Compute the numerator:
    #   count of the n-gram + k
    # Adding k ensures unseen n-grams still get a small nonzero probability
    numerator = ngram_counts.get(ngram, 0) + k

    # Compute the denominator:
    #   count of the context + k * vocab_size
    # The term k * vocab_size redistributes some probability mass
    # across all possible next words in the vocabulary
    denominator = context_counts.get(context, 0) + k * vocab_size

    # Return the smoothed probability of this n-gram
    # Formula: P(w_n | context) = (count(ngram) + k) / (count(context) + k * |V|)
    return numerator / denominator
```

```
Start coding or generate with AI.
```

## 🧩 5. Linear Interpolation of N-gram Orders

### ◆ Theory

High-order N-grams (like trigrams) are powerful but **sparse**.
Low-order N-grams (like unigrams) are reliable but **less specific**.

**Linear interpolation** combines them:

$$P_{interp}(w_i|w_{i-2}, w_{i-1}) = \lambda_3 P(w_i|w_{i-2}, w_{i-1}) + \lambda_2 P(w_i|w_{i-1}) + \lambda_1 P(w_i)$$

where $\lambda_1 + \lambda_2 + \lambda_3 = 1$.

This balances specificity and generality, improving robustness.

### ◆ Function

```
1 from typing import Tuple, Dict
2
3 def interpolated_probability(
4     ngram: Tuple[str, ...],
5     uni_counts: Dict,
6     bi_counts: Dict,
7     tri_counts: Dict,
8     uni_contexts: Dict,
9     bi_contexts: Dict,
10    tri_contexts: Dict,
11    vocab_size: int,
12    lambdas: Tuple[float, float, float] = (0.1, 0.3, 0.6),
13    k: float = 1.0
14 ) -> float:
15    """
16    Computes the final probability of a trigram using interpolation
17    between unigram, bigram, and trigram probabilities.
18    """
19
20    # Extract unigram, bigram, and trigram versions of the input n-gram
21    # Example: if ngram = ("i", "love", "nlp")
22    #    unigram = ("nlp",)
23    #    bigram  = ("love", "nlp")
24    #    trigram = ("i", "love", "nlp")
25    unigram = (ngram[-1],)
26    bigram = ngram[-2:]
27    trigram = ngram
28
29    # Compute smoothed probabilities for each model order
30    # Unigram → single word probability
31    p1 = smoothed_probability(unigram, uni_counts, uni_contexts, vocab_size, k)
32    # Bigram → conditional probability of last word given previous one
33    p2 = smoothed_probability(bigram, bi_counts, bi_contexts, vocab_size, k)
34    # Trigram → conditional probability of last word given two previous ones
35    p3 = smoothed_probability(trigram, tri_counts, tri_contexts, vocab_size, k)
36
37    # Unpack interpolation weights (λ1, λ2, λ3)
38    # These control how much influence each model order has
39    λ1, λ2, λ3 = lambdas
40
41    # Compute the final interpolated probability
42    # Formula: P = λ1 * P_unigram + λ2 * P_bigram + λ3 * P_trigram
43    # Ensures robust probabilities even for unseen higher-order n-grams
44    return λ1 * p1 + λ2 * p2 + λ3 * p3
45
```

## ✳️ 6. Perplexity Evaluation

### ◆ Theory

**Perplexity** measures how "surprised" a model is by new data.

It's the **exponential of the average negative log-likelihood:**

$$PP = \exp\left(-\frac{1}{N}\sum_i \log P(w_i|context)\right)$$

Lower perplexity = better prediction.

It's a standard metric for comparing language models.

### ◆ Function

```python
 1 def calculate_perplexity_interpolated(
 2     test_tokens: List[str],
 3     uni_counts: Dict,
 4     bi_counts: Dict,
 5     tri_counts: Dict,
 6     uni_contexts: Dict,
 7     bi_contexts: Dict,
 8     tri_contexts: Dict,
 9     vocab_size: int,
10     lambdas: Tuple[float, float, float] = (0.1, 0.3, 0.6),
11     k: float = 1.0
12 ) -> float:
13     """
14     Calculates the perplexity of a language model that uses interpolated
15     unigram, bigram, and trigram probabilities.
16     Lower perplexity → better model generalization.
17     """
18
19     # Step 1: Build all trigrams from the test token sequence
20     # Example: ["<s>", "i", "love", "nlp", "</s>"]
21     # → [("<s>", "i", "love"), ("i", "love", "nlp"), ("love", "nlp", "</s>")]
22     test_trigrams = build_ngrams(test_tokens, 3)
23
24     # Step 2: Initialize variable to accumulate the sum of log probabilities
25     log_prob_sum = 0
26
27     # Step 3: Loop through each trigram in the test data
28     for ngram in test_trigrams:
29         # Compute interpolated probability for this trigram
30         prob = interpolated_probability(
31             ngram,
32             uni_counts, bi_counts, tri_counts,
33             uni_contexts, bi_contexts, tri_contexts,
34             vocab_size, lambdas, k
35         )
36
37         # Add the log of the probability to the running total
38         # (Using logs avoids underflow when multiplying many small probabilities)
39         log_prob_sum += math.log(prob)
40
41     # Step 4: Compute total number of trigrams
42     N = len(test_trigrams)
43
44     # Step 5: Compute perplexity using the formula:
45     #   Perplexity = exp(- (1/N) * Σ log(prob))
46     # Lower perplexity → the model is assigning higher probabilities to the correct words
47     return math.exp(-log_prob_sum / N)
48
49
```

## 🧩 7. Text Generation

### ◆ Theory

Once the model knows how likely words follow each other, we can **sample** new text.

Starting with a context ( `<s> <s>` for trigrams), we repeatedly draw the next word based on its probability until we reach `</s>` .

This process mimics natural language creation.

### ◆ Function

```python
1 from typing import List, Tuple, Dict
2 import random
3
4 def generate_text_interpolated(
5     vocab: List[str],
6     uni_counts: Dict,
7     bi_counts: Dict,
8     tri_counts: Dict,
9     uni_contexts: Dict,
10     bi_contexts: Dict,
11     tri_contexts: Dict,
12     length: int = 20,
13     lambdas: Tuple[float, float, float] = (0.1, 0.3, 0.6),
14     k: float = 1.0
15 ) -> str:
16     """
17     Generates text from a trigram language model using interpolated probabilities.
18     """
19
20     # Step 1: Initialize the starting context with two sentence-start tokens
21     # These act as initial history for trigram generation
22     context = ["<s>", "<s>"]
23
24     # Step 2: Initialize an empty list to hold the generated words
25     output = []
26
27     # Step 3: Loop until reaching the desired number of words (length)
28     for _ in range(length):
29         probs = []  # list to store (word, probability) pairs
30
31         # Step 4: For each word in the vocabulary, compute its probability
32         for word in vocab:
33             # Construct the trigram: last two context words + candidate next word
34             ngram = tuple(context[-2:] + [word])
35
36             # Get the interpolated probability for this trigram
37             p = interpolated_probability(
38                 ngram,
39                 uni_counts, bi_counts, tri_counts,
40                 uni_contexts, bi_contexts, tri_contexts,
41                 len(vocab), lambdas, k
42             )
43
44             # Store the word and its probability
45             probs.append((word, p))
46
47         # Step 5: Normalize probabilities so they sum to 1
48         total = sum(p for _, p in probs)
49         probs = [(w, p / total) for w, p in probs]
50
51         # Step 6: Randomly sample the next word, weighted by probabilities
52         next_word = random.choices(
53             [w for w, _ in probs],       # list of candidate words
54             weights=[p for _, p in probs]  # their corresponding probabilities
55         )[0]
56
```

```
57          # Step 7: If the end-of-sentence token appears, stop generation
58          if next_word == "</s>":
59              break
60
61          # Step 8: Add the chosen word to the output and update context
62          output.append(next_word)
63          context.append(next_word)
64
65      # Step 9: Join generated words into a single string and return it
66      return ' '.join(output)
67
```

## ❖ 8. Main Script

### ◆ Theory

This brings all components together:

1. Load data and tokenize it.

2. Train unigram, bigram, trigram models.

3. Evaluate performance using **perplexity**.

4. Generate example text using learned probabilities.

### ◆ Function

```
1 !wget https://raw.githubusercontent.com/karpathy/char-rnn/master/data/tinyshakespeare/input.txt -O sample_text.txt
2
```

```
--2025-11-12 07:34:35--  https://raw.githubusercontent.com/karpathy/char-rnn/master/data/tinyshakespeare/input.txt
Resolving raw.githubusercontent.com (raw.githubusercontent.com)... 185.199.108.133, 185.199.109.133, 185.199.110.133, ...
Connecting to raw.githubusercontent.com (raw.githubusercontent.com)|185.199.108.133|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 1115394 (1.1M) [text/plain]
Saving to: 'sample_text.txt'

sample_text.txt     100%[===================>]   1.06M  --.-KB/s     in 0.006s

2025-11-12 07:34:35 (180 MB/s) - 'sample_text.txt' saved [1115394/1115394]
```

```
1 def main():
2     """
3     Main function:
4     ----------------
5     Coordinates the full language modeling pipeline:
6     1. Loads and preprocesses text
7     2. Trains unigram, bigram, and trigram models
8     3. Evaluates models using interpolated perplexity
9     4. Generates new sample text using interpolation
10
11    Acts as a demonstration and testing harness for the n-gram model functions.
12    """
13
14    # ------------------------------------------------------------
15    # 1. Load and preprocess text
16    # ------------------------------------------------------------
17    with open("/content/sample_text.txt", "r", encoding="utf-8") as f:
18        text = f.read()
19
20    # Tokenize the input text (split into words or symbols)
21    tokens = tokenize_text(text)
22
23    # Build the vocabulary and calculate its size
24    vocab = list(set(tokens))
25    vocab_size = len(vocab)
26
27    # Split tokens into training (80%) and testing (20%) sets
28    split = int(0.8 * len(tokens))
29    train_tokens, test_tokens = tokens[:split], tokens[split:]
```

```
30
31    # --------------------------------------------------------------
32    # 2. Train N-gram models
33    # --------------------------------------------------------------
34    # Train unigram, bigram, and trigram models on the training tokens
35    uni_counts, uni_ctx = train_ngram_model(train_tokens, 1)
36    bi_counts, bi_ctx = train_ngram_model(train_tokens, 2)
37    tri_counts, tri_ctx = train_ngram_model(train_tokens, 3)
38
39    # --------------------------------------------------------------
40    # 3. Evaluate model performance using perplexity
41    # --------------------------------------------------------------
42    # Use interpolated model (combining unigram + bigram + trigram probabilities)
43    # Lambdas are the weights for each n-gram level
44    perplexity = calculate_perplexity_interpolated(
45        test_tokens,
46        uni_counts, bi_counts, tri_counts,
47        uni_ctx, bi_ctx, tri_ctx,
48        vocab_size,
49        lambdas=(0.1, 0.3, 0.6)  # weights for 1-gram, 2-gram, 3-gram
50    )
51
52    print(f"\nPerplexity: {perplexity:.2f}")
53
54    # --------------------------------------------------------------
55    # 4. Generate new text using interpolated n-gram model
56    # --------------------------------------------------------------
57    print("\nGenerated Text:")
58    generated = generate_text_interpolated(
59        vocab,
60        uni_counts, bi_counts, tri_counts,
61        uni_ctx, bi_ctx, tri_ctx,
62        length=30,              # number of tokens to generate
63        lambdas=(0.1, 0.3, 0.6)
64    )
65    print(generated)
66
```