# An Approach to Modelling English Sentence Structures Based on the Combination of Graph Theory and Probability

**¹Ambreen Zehra Rizvi, ²Nazra Zahid Shaikh**

| Article Details | ABSTRACT |
|---|---|

**Ambreen Zehra Rizvi**
Assistant Professor, Faculty of Engineering, Science & Technology, Hamdard University Main Campus, Karachi, Pakistan.
ambreen.zehra@hamdard.edu.pk

**Nazra Zahid Shaikh**
Senior Lecturer, Department of English, Faculty of Social Sciences and Humanities, Hamdard University Main Campus, Karachi, Pakistan.
nazra.zahid@hamdard.edu.pk

This paper studies the uses of graph theory and probability theory to model English sentence structures. Through a representation of sentences as directed graphs (parsing into dependency trees) and analysis of syntactic patterns using probability models, we also illustrate, how mathematical parametric modelling can increase our understanding of linguistic data. We take a look at the dependency grammar, as well as at the Markov models and the probabilistic contextfree grammars (PCFG) to investigate how mathematical tools can forecast and explain sentence construction. These results indicate that graph-theoretical and probabilistic tools are very strong tools for NLP, syntax-parsing and machine-translation.

## INTRODUCTION

There is an interesting relationship between linguistics and mathematics when it comes to language structures. In contrast to linguistics, which has traditionally interested itself with the rules of and patterns in human language, mathematics offers formalisms for modeling, analyzing, and predicting elements of language. Syntax, one of the subfields of linguistic studies,- the study of sentence structure- shares certain properties with other branches of science, mainly the mathematical ones, in the sense that it is based on hierarchies and recursivity.

As with most languages, English follows a specific set of rules when making sentences with words. These rules can be studied systematically with the aid of two main mathematical methods:

1. **Graph Theory** – This sub-branch of Discrete mathematics is very good for representing structural relationship between words in a sentence. If sentences are modeled as directed graphs (e.g. dependency trees or constituency parse trees), linguists and computer scientists are able to understand syntactic hierarchies, identify grammatical errors and design natural language parsing algorithms.

2. **Probability Theory** – Because language is uncertain (words and phrases have probabilistic behavior instead of crisp deterministic behavior), probabilistic models like Markov chains, n-gram models, and probabilistic context-free grammars (PCFGs) are used to predict probable sequences of words, to disambiguate sentence structure, and to enhance machinedriven language processing.

The inter-weaving of these mathematical approaches has disrupted disciplines such as the computational linguistics, the natural language processing (NLP), and the AI based text generation. For example, current language models such as GPT-4 and BERT exploit uncertain graph-based technology to understand and produce human-like text.

This paper explores how both graph theory and probability theory may be used to model the structure of English sentences, studying:

· Syntax and its representation in terms of dependency and constituency graphs.

· The contribution of statistical language models to grammaticality prediction.

· The practical uses of such models within strong parsing systems, MT systems and AI linguistics.

This work demonstrates the way in which the integration of linguistic theory with mathematical formalism highlights the facility with which language can be better understood when reformulated in quantitative terms, providing new apparatus for theoretical analysis and for practical applications to artificial intelligence.

This paper investigates:

· The graph-based method of dependency trees for describing sentence structure.

· How probability models (n-grams, PCFG) are used to predict grammatical sentences.

· Computational linguistics and artificial intelligence (AI) applications.

## GRAPH THEORY TO MODEL SYNTAX

The Definition of Dependency Grammar as Dirrected Graphs

According to the dependency grammar, sentences are DAGs: each word in a sentence forms a node to be connected with other nodes by directed edges, indicating the dependencies between words. While constituency-based theories like the constituency grammar group words into nested phrases (e.g., noun phrases, verb phrases), dependency grammar instead uses binary word-to-word relations (ones word's role in the sentence is being a parent of another word), making relations more direct and thus easier to determine.

Structural Representation

**Nodes:** Every word in the sentence is a node in the graph.

For example, in the sentence "The cat chased the mouse", the nodes are: "the, cat, chased, the, mouse".

**Edges:** Directed links representing syntactic connections (generally from a head word to a dependent).

Example:

"cat" → "the" (relationship between a determiner and a noun)

"chased" → "cat" (subject-verb)31 in Fig.

"chased" → "mouse" (subject-verb dependency)

Example: Dependency Tree of "The cat chased the mouse."

```
              chased
              /    \
           cat mouse
           /    \
          the   the
```

In the first sentence "chased" is the root of the sentence (main verb) and governs two terms: Subject("cat") and Object("mouse"), here are the two phrases with name of the constituents (NP or N'P) inside the brackets. The determiners ("the") belong to the nouns.

Advantages Relative to Constituency Grammar

**Simple:** Dispenses with intermediate phrases so simpler.

**Direct Associations**: This makes explicit which words modify or head others (e.g., which adjectives are directly connected to which nouns).

**Computational Need:** These algorithms can parse the dependency trees in linear time such as Eisner algorithm for the NLP.

## FORMAL PROPERTIES

**Crossing:** Arcs do not cross for most of the well-formed sentences (projective). Crossed arcs (non-projective trees) are often associated with complex Phenomena like dislocations or free word order (in German, Latin, etc.).

## TREE CONSTRAINTS

· Single root (usually the main verb).

· No cycles (i.e., ensuring grammatical hierarchy).

· Linked structure (words must all be connected).

Computational Linguistic Applications

**Dependency Parsing:** MaltParser and spaCy use graph-based algorithms to parse sentences.

**Machine Translation:** Dependencies trees are helpful for aligning sentences across languages by tracking the syntactic roles.

**Grammar Check(ing):** Identifies errors such as misagreement of subject and verb by examining edge connections.

The representation of syntax as a graph is such that linguists and computer scientists can take advantage of graph algorithms (e.g., topological sort, shortest path) to better model natural language than the with rule-based systems.

Graph-Based Syntactic Structures in Syntactic Modeling Applications

Graph theory offers a variety of powerful tools to analyze and manipulate syntactic structures in natural language. Sentences, when viewed as directed graphs, enable graph-theoretic concepts for addressing critical problems in syntax analysis for linguists and computational scientists. Here, we focus on three key applications:

## CYCLE DETECTION FOR VALIDATING GRAMMATICALITY

One of the core limits of dependency grammar is that syntactic structures must be acyclic graphs – that is, a word cannot depend on itself (either directly or indirectly). Graph theory provides fast algorithms to detect cycles, guaranteeing the sentences respect the grammar rules. Example of an Illegal Cycle:

·"The cat that chased the mouse ran."

· If this is naively parsed, you have a circular dependency, with "cat" depending on "chased", which depends on "mouse" and so being already assigned in parallel to both "chased" and "mouse".

## ALGORITHMIC SOLUTION

· Cycles in DFS tree can be found in $O(V+E)$ time.

· Package tools such as Stanford Parser and spaCy reject cyclic dependencies in parsing automatically.

It is an important application for grammar checking, sentence correction systems and maintaining well-formed syntactic trees in NLP pipelines.

## TREE WIDTH MEASURES FOR THE ANALYSIS OF SYNTACTIC COMPLEXITY

Graph-theoretic measures can also be used to quantify the structural complexity of a sentence, including, for example:

**Tree Depth:** The maximum distance from root (main verb) to leaf node (e.g., det or adj).

Example: "The quick brown fox jumps over the lazy dog." has more depth than "The cat sleeps."

**Tree Width:** Measures the "branchiness" of a sentence ( for example multiple modifiers for one noun ej. increase the width).

This has also application in psycholinguistics to investigate cognitive load in sentence comprehension.

**Crossing Edges:** Non-projective sentences are more complex (dependencies cross) and are not

so rare in languages with free word order (e.g., German).

These metrics help in:

- Text Complexity Analysis (eg Simplification of medical or legal text).

- Language acquisition studies (child versus adult sentence patterns).

- Graph-Based Methods for Machine Parsing

Converting plain text to dependency trees is computationally intensive, so high-performance algorithms are needed. Two key approaches are:

## EISNER'S ALGORITHM

A dynamic programming algorithm that decodes sentences in $O(n)$ time.

Ensures projective (non-crossing) trees, which is very good for English.

MST Parsing from maximum spanning tree (MST):

Considers dependency parsing as a graph optimization task.

Weights edges (e.g., using machine learning) and finds the tree with the highest score.

Real-World Use Cases:

Google's SyntaxNet feels the sentence using graph-based parsing.

Voice Assistants (e.g., Siri, Alexa) use these parsers to understand user commands.

## ADDITIONAL APPLICATIONS

**Coreference Resolution:** Graph clustering methods connect pronouns with corresponding referents (e.g., "She" → "Mary").

**Semantic Role Labeling:** Fills "in the blank _ who did what to whom _ by determining relationships among words.

The utility of graph theory in syntax – whether in cycle detection or parsing - shows that it is an indispensable tool in computational linguistics and NLP. By graphifying sentences, linguists can now take advantage of fifty years of graph analysis algorithms to model, analyze, and process (human) language in a way not previously possible.

## PROBABILITY IN SENTENCE PREDICTION

## MARKOV MODELS ON WORD SEQUENCE

In natural language, probability theory has a solid role to play in the modelling and prediction of sentence structures. The most basic model in this category is the Markov model, under the assumption that the probability of a word occurring is based only on the words previous to it. This is particularly beneficial to modeling local syntactic pattern and sequential sequence of

language.

## MATHEMATICAL BACKGROUND OF N-GRAM MODELS

The most basic Markov models for language are the so-called n-gram model, where the probability of a word sequence is estimated as a function of the previous (n-1) words:

--Unigram Model (1-gram):

P(w_n)

Treats each word separately, with no regard to context

Easiest probability model for word counts

--Bigram Model (2-gram):

P(w_n | w_{n-1})

Predicts current word from last word immediately to the left

Example:

P("the"|"chased") > P("a"|"chased") in normal English

Learns simple local word dependencies

--Trigram Model (3-gram):

P(w_n | w_{n-2}, w_{n-1})

Takes into account 2 previous words while predicting

It is a little bit less bad at capturing long-range dependencies

– More contextual than bigram but sparser data

--Probability Estimation:

(Maximum Likelihood Estimation (MLE) formula:)

$$P(w\n|w\backslash\backslash\{n\text{-}k\backslash\},\backslash\backslash ldots,w\backslash 1) = Count(w\backslash\backslash\{n\text{-}k\backslash\},\backslash\backslash ldots,w\n)/Count(w\backslash\backslash\{n\text{-}k\backslash\},\backslash\backslash ldots,w\backslash\_\backslash\{n\text{-}1\backslash\})$$

Where:

Count() = frequency in training corpus

  - k = n-gram order minus 1 (1 for bigram, 2 for trigram)

  - Numerator = count of full n-gram sequence

  - Denominator = count of context words only **Properties of n-Gram Models**

## UNIGRAM MODELS: SIMPLICITY VS. CONTEXT

Unigram models: Most rudimentary probabilistic language models that generate the probability of each word independently of the rest based on the final output. This model estimates the probability of a word occurring in a text based only on its frequency in the

training corpus and ignores the words that come before or after. Although unigram models are computationally attractive and simple to implement, they are severely limited in their ability to model linguistic structure. Since they cannot capture context between words, they are not applicable to tasks which require syntactic or semantic understanding (e.g., coherent text generation or accurate sentence completion). But unigram models are still useful in simple tasks such as baseline text categorization or as building block of a more sophisticated model.

## BIGRAM MODELS: LOCAL DEPENDENCE

Bigram methods present a coarse modification over unigram methods by incorporating information about immediate neighbors – the pairs of words in the training data. These learn to estimate the probability of a word (a unit of language in written form or speech) based on its second most recent predecessor, and thus capture simple local patterns of language use such as natural adjective-noun combinations or subject- verb aggreement. Although bigrams achieve much better performance than unigrams in several tasks including data prediction, their length-1 lookback window is still restricted and does not capture the long-range dependence of language. For example, they can lack adequate agreement between the subject and verb, as in through intervening phrases and clauses. Nevertheless, bigram models still offer good computational performance and compromise for plenty of simple language processing tasks.

## TRIGRAM MODEL: BEYOND BIGRAM MODEL WITH SPARSITY ISSUES

Trigram model extend the contextual information further by taking into account three adjacent words, which account for more fine-grained linguistic cues than bigram model. This extra context serves to let the trigram models treat some grammatical forms and frequent three-word expressions in the natural language in a more effective manner. But as is always the case with expanded context, the gains don't come without some pain. For larger n-gram orders, the model is faced with severe data sparseness problems because a lot of three-word combinations happen either rarely or not at all in training corpora. This problem of sparseness, will however, make estimate of probabilities of rare or unseen trigrams untrustworthy, and may degrade the performance of the model. Several smoothing methods have been proposed to reduce these disadvantages while retaining the advantages of the extra context.

## MAXIMUM LIKELIHOOD ESTIMATION: PROS AND CONS

The prevailing n-gram probability estimation method, Maximum Likelihood Estimation (MLE), assigns probabilities by taking a count of the n-gram and dividing it by the count of its prefix in the training data. Although MLE yields probabilities that are transparent and

intuitive, there are several important drawbacks of the procedure. Most critically, it gives zero weight to any n-gram that never occurred in the training corpus, which is especially damaging when considering ngrams of high order where the number of possible ensembles is astronomical. In addition Note that MLE overestimates the probability of the observed events while underestimate the probability of the unobserved but possible sequences. Due to these deficiencies, a variety of smoothing techniques are required to develop more reliable probability estimates.

## METHODS OF SMOOTHING AND ROBUST ESTIMATION

## LAPLACE (ADD-ONE) SMOOTHING

The simplest and initial smoothing method, called Laplace smoothing, solves the zero-probability problem adding one to the counts of every possible n-gram, including the one which never shows up in the training data. Although this approach outputs no probability estimate of zero, it tends to overly simplify the probability distribution by putting excessively high mass on unseen events.

This is a very large source of bias for the estimates of the model (especially if the size of the vocabulary is high, where counting possible n-grams becomes tedious). Given its potential limitations, Laplace smoothing does, however, provide a useful baseline method and is computationally tractable for small-scale problems.

## GOOD-TURING DISCOUNTING

Good-Turing smoothing offers a more elegant solution, redistributing probability mass from frequent to rare and unseen events. The method estimates the probability of unseen n-grams by counting the frequency with which n-grams of count one appear in the training-text. A better approximati on of the true distribution of language is obtained by this more balanced distribution of probability: i n reality, some unobserved sequences are arguably plausible. The Good-Turing performs well especially over small corpora and has been reported to give reliable estimates when properly applied. Nevertheless, it depends on the treatment of the count statistics, and it is a computationally demanding solution for high n-gram orders.

## BACKOFF VS INTERPOLATIONS

Backoff and interpolation techniques provide smoother treatment of sparse data. Backoff methods fall back to lower-order n-gram counts if higher-order counts are missing or become zero and interpolation methods interpolate estimates of all n-gram orders, up to the current one. A popular realization known as Katz backoff employs a more intricate discounting mechanism

and backs off the lower-order models. Interpolation methods like Jelinek-Mercer smoothing weight the contribution of varying order n-grams according to their trustworthiness. Such methods tend to surpass the simple smoothing techniques since they make the utilization of the available data more economical, and they generate more accurate predictions on probability of different contexts.

## KNESER-NEY SMOOTHING

Kneser-Ney smoothing, one of the most successful smoothing methods for n-gram models, provides a new, distribution-based way for estimating lower-order models. Instead of basing backoff probabilities on frequencies of words, it bases them on the number of different contexts a word has appeared in, which seems to model the true distribution of words in language better. Such an approach is particularly useful to deal with stop words which occur in fewer contexts and rare words which occur in more contexts. Modified Kneser-Ney smoothing techniques, including interpolated smoothing, provide further improvement by employing absolute discounting but in combination with more sophisticated backoff model. These more sophisticated smoothing techniques have become ubiquitous in modern n-gram language models, and they greatly enhance the predictive accuracy and resilience of the statistic.

All of these smoothing approaches deal with different facets of the problem of data sparsity, where more sophisticated methods usually achieve better performance at the expense of being timeconsuming. The selection of the smoothing method is based on criteria such as the size of the training corpus, the order of the n-gram model and the particular needs of the application. Stateof-the-art language processing systems frequently use a combination of smoothing methods, or incorporate them along with more sophisticated neural language models to ensure best task performances.

A key problem in n-gram modeling is how to deal with so called sparse data - that is, word ngrams that have zero count in the training corpus. There are several approaches to smoothing:

## LAPLACE (ADD-ONE) SMOOTHING:

Adds 1 to every count (no zero probabilities).

## GOOD-TURING DISCOUNTING

Moves probability mass 'downhill' from common to infrequent events.

## BACKOFF AND INTERPOLATION

Utilizes lower-order n- grams when higher-order counts are zero. applications in language and

speech technology

Applications Markov models provide theoretical underpinning of many practical applications:

Autocorrect Systems:

These systems are able to identify and correct unlikely combinations (e.g. predict the' after chased' instead of `a') by computing the probability of sequences of words.

## SPEECH RECOGNITION

Hidden Markov Models (HMMs) decode the spoken word by determining the most likely sequence of words from an acoustic input.

## TEXT GENERATION

Simple forms of language models can issue fairly decent sounding sentences by sampling from ngram distributions for example.

## MACHINE TRANSLATION

Assists in evaluating fluency of translated output by scoring the likelihood of target language word sequences.

## LIMITATIONS AND EXTENSIONS

Traditional n-gram-based models are effective, but they have a number of limitations:

· Long-range Dependencies:

Inability to capture dependencies beyond the n-gram (e.g. subject verb agreement over an extended distance).

· Context Insensitivity:

Effects exactly the same changes in all occurrences of one vaccine name, irrespective of the semantic or contextual significance of the string.

These are now dealt with in contemporary scripts by:

· Neural Language Models:

Apply recurrent or transformer models to represent longer range patterns.

· Conditional Random Fields:

Use all sorts of other linguistic information also.

This is a moment of history in the evolution of this probabilistic framework, for the current trend is on combining neural approaches with classic Markovian assumptions in order to do better language modeling.

## PCFGS (PROBABILISTIC CONTEXT-FREE GRAMMARS)

Probabilistic context-free grammars are a major stepforward in syntactic model- ing as they extend standard context-free grammars with probabilistic information. The extension enables an analysis of sentence structures at the level of quantities and is an important disambiguation tool for natural language processing.

## THEORETICAL FOUNDATIONS

PCFGs work under the same principles as regular CFGs while adding a probability distribution over the production rules. Each grammar rule has an associated probability value, resulting in a stochastic grammar in which:

· the probability of a parse tree P is the product of the rule probabilities used in its derivation

· WordsAll the words used must be terminal words. For example, there is no word "John sleeps" in the grammar above, only phrases with structure and meaning PKKU309 Natural Language Processing The probability of the sentence is the sum of the probabilities of all possible parse trees that generate it.

This probabilistic model enables the grammar to:

Ordering alternative parses of a sentence with ambiguous structure

Learn a model of the relative frequency of syntactic structures

Quantify the parse quality.

## RULE PROBABILITY ESTIMATION

The production rule probabilities are often learned from a parsed corpora (treebank) using maximum likelihood estimation. For a rule A → α:

$$P(A \to \alpha) = \frac{count(A \to \alpha)}{count(A)}$$

Where:

· count(A → α) is the count of the rule in the training data

· count(A) is the number of complete expansions of non-terminal A.

For thith ro aore frith acthish of roathiatha usath mory roramate sporacth in use hith atan language.yielding higher probabilities to more common syntactic patterns seen in genuine language usage.

## PARSING ALGORITHMS FOR PCFGS

There are two major methods for parsing with PCFGs:

Cocke-Kasami-Younger (CKY) Algorithm: In contrast, the corresponding CKY algorithm is given by.

· Bottom-up DP approach

· Build a parse chart from all possible constituents

· Re-estimates the Probability in parsing-sectional incremental manner.

· runs in $O(n^3)$ time w.r.t. sentence length n

## EARLEY'S ALGORITHM

· A top-down with bottom-up filtering.

· It is more effective for some grammar types

· Contains probability as well

## HANDLING STRUCTURAL AMBIGUITY

Probabilistic context-free grammars (PCFGs) are good for resolving such syntactic ambiguities by comparing probabilities. Consider the sentence:

"Police help dog bite victim"

## POSSIBLE INTERPRETATIONS:

[Police Assisting [victim of dog bite]] (more likely)

[Police [help dog] [bite victim]] (less likely)

It is those probabilities assigned to the VP → V NP and VP → V NP NP rules which determine the preferred parse.

## LEXICALIZED EXTENSIONS

Simple PCFGs were also improved by means of lexicalization:

· Each nonterminal has a head word

· Lexical Item: Rule probabilities are conditional on lexical items.

· Encodes conceptual relations in the representations of the syntactic structures

For example:

VP(vb:help) → V(help) NP(nn:dog) NP(nn:victim)

## APPLICATIONS IN MODERN NLP

PCFGs still have significant applications in:

### GRAMMAR DEVELOPMENT

· Construction of probabilistic grammars for new languages Parser Construction:

· Stanford Parser and other machine learning based parsers.

### LANGUAGE MODELING

· Syntactic information in language models

### GRAMMAR CHECKING

· Recognizing unlikely parses as potential errors

### EVALUATION AND LIMITATIONS

PCFGs are powerful but suffer from several issues:

### CONTEXT INDEPENDENCE

· Rule probabilities are not aware of wider context

### LEXICAL GAPS

· Combat inappropriate word groupings behind the scenes Structural Sparsity:

· Some rare constructions may not be well accounted for

Recent work tries to mitigate these drawbacks by using PCFGs and neural network combined while still keeping the grammatical framework interpretable.

### NATURAL LANGUAGE PROCESSING FOR MIXED-INITIATIVE SYSTEMS

Today's NLP systems embody a rich combination of graph-theoretic and probabilistic methods, which yield hybrid architectures resulting in state-of-the-art capabilities in language processing and generation. These systems integrate:

### ATTENTION-BASED GRAPH STRUCTURED NETWORKS

AweSome: These days, model architectures that are based on transformers and those transformer architectures that include gaph-like attention (such as BERT and GPT-4) owe to the above in multiple significant ways:

### DEPENDENCY-AWARE ATTENTION

• Self-attention layers automatically capture dependency relationships between words · The attention heads have different syntactic relations covered (subject-verb, modifier-head)

· Empirical: some heads match traditional dependency grammars spectacularly.

### IMPROVING GRAPH NEURAL NETWORK

· In some architectures dependency parse is directly injected.

· Graph convolutional layers work over parsers with attention on top. · Performance increases

on tasks that involve structural understanding

## RELATIONAL BIAS IN ATTENTION

· Positional encodings combined with syntactic distance features

· Dependency path limited attention masks

· Model is given the ability to better model long distance dependencies

## NEURAL PROBABILISTIC LANGUAGE MODELS

The probabilistic underpinnings of modern NLP systems are realised in:

## DISTRIBUTIONS OF CONTEXTUAL WORDS

· Models out probability distributions over vocab.

· Output Diversity is controlled by Temperature Sampling

· Beam search keeps multiple high-probability sequences

## UNCERTAINTY QUANTIFICATION

· Scores for confidence of predictions

· Output distribution entropy measures · Calibration methods for realistic probabilities

## STRUCTURED PREDICTION

· Models of joint probabilities for sequences

· Conditional random fields for structured prediction

· Viterbi-like decoding methods

Architectural Integration

These are connected through:

Frameworks for Multi-Task Learning: Multi-task learning gains popularity in recent years. · Syntax and semantics objectives joint training

· Multi-task shared representations

· Gradient balancing now!GenerationStrategy (up to second order)!

## HYBRID PARSING SYSTEMS

· Traditional parsers conditioned on neural network feature extractors

· Neural reranking using graph-based representations

· Methods combining different approaches together into an ensemble.

## EMERGENT SYNTACTIC LEARNING

· Attention patterns analysis and learned grammar identification

· Investigations suggesting hierarchical processing

· Continuous representations of syntactical rules

Applications and Implications This integration enables:

## LESS BAD MACHINE TRANSLATION

· Reordering is improved by attention that is aware of the syntax

· Fluency is preserved by probabilistic decoding

## BETTER QUESTION ANSWERING

· Reasoning over text using graphs

· Oracle-based confidence scoring for answers

## CONTROLLED TEXT GENERATION

· Syntax-constrained decoding

·Risk-averse Style Transfer

The integration of graph-based structural awareness with neural probabilistic modeling is currently the most advanced approach, providing the interpretability of formal linguistic analysis and statistical flexibility simultaneously. This parallel effort remains the main engine of progress for natural language processing in all its aspects.

## JOINT ADVANCES IN NATURAL LANGUAGE PROCESSING

Contemporary NLP systems have obtained impressive results by leveraging graph-based structural analysis and neural probabilistic modeling. This combination of both paradigms combines the advantage of both in an effort to develop a stronger and more accurate language processing system. The heart of such models like BERT and GPT-4' is a clever amalgam of these strategies which is able to capture not only syntactic relations but also semantic probabilities.

The attention mechanism utilized in the Transformer models inherently captures the dependencies among words and forms a dynamic graph which varies with the context. Several attention heads in these models focus on (parts of) different syntactic relations, and little is known about which heads are dedicated to certain syntactic relations: e.g., it has been empirically demonstrated that some heads nearly always capture relations, such as subject-verb relationships, and others concentrate more on modifier-head relations. This emergent behavior shows how neural architectures are able to implicitly learn grammatical structure without supervision, though some systems augment this with direct use of dependency parse

information.

On the probabalistic side, recent language models write text by working with probability distributions over vocabulary elements. They ensure that uncertainty has been estimated and maintained throughout the generation process (e.g.: leveraging techniques such as temperature scaling or beam search that allow to flexibly trade off between creativity and coherence). The probability-based formulation of this method makes it possible to offer confidences calibrated for applications such as question answering and to provide controlled generation settings by tuning the output distributions. The models' command of what word sequences is likely also comes from being exposed to large corpora of text, from which they learn the statistical tendencies of the language that are expressed therein.

State-of-the-art systems combine these by embedding them in multi-task learning frameworks, which co-optimise syntactic and semantic objectives. Some networks use explicit graph neural network components in addition to the classical attention mechanisms in the form of hybrid models capable of reasoning not only about local word probabilities, but also about global sentence structure. Others apply probabilistic graphical models to reranking or constraint lower neural network outputs, combining traditional NLP with modern deep-learning methods.

This union of methodologies has also resulted in performance gains in various NLP tasks. In the machine translation, syntax-aware attention mechanisms facilitate word reordering among language-pairs and the model performs probabilistic decoding that preserves fluent naturalness. For question answering systems, this combination allows for an improved extraction of relationships and the estimation of higher confidences. The field is still developing as researchers find new ways to blend structural and probabilistic information, and push the limits of what is possible in understanding and generating language.

## CONCLUSION

When we fuse graph theory and probability theory something amazing happens: for the first time in the history of human thought we posses the capability to model and analyze English syntax with some degree of computational accuracy! These formalisms offer complementary benefits that cover different aspects of linguistic structure and behavior. Graph-based methods, and specifically the dependency grammar representations, provide an intuitive and computationally efficient way to represent the hierarchical relationship of words in a sentence. By modeling sentences as directed graphs, graph algorithms known to linguists and computer

scientists can now be used to effectively parse, detect cycles in, and analyze the complexity of sentence structure, making them feasible for rigorous syntactic processing in NLP.

This stochastic aspect provides an important element of natural language use for syntactic modeling. Probability distributions enable systems to measure the relative likelihood of different grammatical solutions, disambiguate among competing parses, and produce fluent, context appropriate language. We all know that statistical patterns inferred from large corpora can heavily improve the performance of language processing systems, as shown by Markov models and probabilistic context-free grammars. This kind of probabilistic model is especially useful since it reflects the kind of variability and uncertainty that exist in natural human language, where there are often several equally correct syntactic constructions that coexist with different probabilities of usage.

The most recent big leaps in this field have been obtained by mixing these techniques with solutions and methods from the deep learning scene. The recent transformer-based models (e.g. BERT, GPT) encode both structural and statistical information in their neural weights and in the attention mechanism. These models show that the future of parsing is probably in hybrid systems which have the interpretability of graph-based systems and the pattern matching of neural networks. Possible directions for future research could include developing more direct methods for integrating syntactic graphs into neural architectures, finding improved techniques for probabilistic calibration of neural predictions, and exploring the extent to which these formal models can provide explanations for the syntactic knowledge learned by large language models.

Combining linguistics, math and computer science in this way the interdisciplinary approach continues to extend the scope for what is feasible in the area of language technology. The more advanced we are able to make our representations and algorithms for handling this kind of knowledge, the closer we are to building artificial intelligence systems capable of understanding and producing human language with nearly human-level competence. The mathematical frameworks of graph theory and probability, which have already been indispensable for the investigation of communicative syntax, will no doubt remain as central tools – both in the practical problem of implementation of the ideas presented here and in the theoretical analysis of the multifaceted phenomenon that is human syntax.

## REFERENCES

Jurafsky, D., & Martin, J. H. (2023). Speech and language processing (3rd ed.). Pearson.

Manning, C. D., & Schütze, H. (1999). Foundations of statistical natural language processing. MIT Press.

Klein, D., & Manning, C. D. (2003). Accurate unlexicalized parsing. Proceedings of the 41st Annual Meeting of the Association for Computational Linguistics (ACL), 423–430. https://doi.org/10.3115/1075096.1075150

Eisner, J. (1996). Three new probabilistic models for dependency parsing: An exploration.

Proceedings of the 16th International Conference on Computational Linguistics (COLING), 340– 345. https://doi.org/10.3115/992628.992688

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., & Polosukhin, I. (2017). Attention is all you need. Advances in Neural Information Processing Systems (NeurIPS), 30, 5998–6008. https://doi.org/10.48550/arXiv.1706.03762

Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (NAACL), 4171–4186. https://doi.org/10.18653/v1/N19-1423

Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., ... & Amodei, D. (2020). Language models are few-shot learners. Advances in Neural Information Processing Systems (NeurIPS), 33, 1877–1901. https://doi.org/10.48550/arXiv.2005.14165