

# PatMaN: A Technical Guide to a High-Throughput Short Sequence Search Tool

**Author:** BenchChem Technical Support Team. **Date:** December 2025

## Compound of Interest

Compound Name: Patman

Cat. No.: B1221989

[Get Quote](#)

For Researchers, Scientists, and Drug Development Professionals

This in-depth technical guide provides a comprehensive overview of **PatMaN** (Pattern Matching in Nucleotide databases), a powerful command-line tool for rapid alignment of large sets of short nucleotide sequences against extensive databases, such as whole genomes. This document details the core algorithm, experimental protocols, performance metrics, and provides visualizations of its operational workflow and underlying logic. The C++ source code for **PatMaN** is available under the GNU General Public License.[\[1\]](#)[\[2\]](#)

## Core Concepts

**PatMaN** is designed for efficiency when searching for numerous short nucleotide sequences, accommodating a predefined number of mismatches and gaps.[\[1\]](#)[\[2\]](#)[\[3\]](#) It is particularly well-suited for applications such as microarray probe mapping, transcription factor binding site identification, and miRNA target analysis. The tool reads both query and database sequences in FASTA format and outputs the alignments in a tab-separated format.

The core of **PatMaN**'s functionality lies in its implementation of a non-deterministic automata matching algorithm built upon a keyword tree of the search strings. This approach allows for exhaustive searches without the heuristic limitations of seed-based alignment methods, which can be crucial when dealing with very short sequences or when alignments with mismatches or gaps are expected.

# The PatMaN Algorithm

**PatMaN**'s algorithm can be broken down into two main phases: keyword tree construction and database searching.

## 2.1. Keyword Tree Construction

Initially, **PatMaN** constructs a keyword tree from the provided set of short query sequences. Each path from the root to a leaf in this tree represents a unique query sequence. To account for searches on both strands of a DNA database, the reverse complement of each query sequence is also added to the tree.

If the user enables the ambiguity flag, the tree is expanded to include all possible nucleotide bases at ambiguous positions within the query sequences. Otherwise, only the standard IUPAC ambiguity code 'N' is recognized and is treated as a mismatch.

## 2.2. Database Searching with a Non-Deterministic Automaton

Once the keyword tree is built, **PatMaN** processes the target database sequence one base at a time. It maintains a list of partial matches, each represented by a node in the keyword tree and an associated edit distance (the number of mismatches and gaps).

For each base in the database sequence, the algorithm attempts to extend all current partial matches by traversing the corresponding edge in the keyword tree. If a perfect match occurs, the partial match is advanced to the next node with no change in the edit distance. In the case of a mismatch or a gap, a new partial match is created with an incremented edit distance, as long as this distance remains below the user-defined threshold. This process allows for the simultaneous exploration of all possible alignments for all query sequences at each position in the database.

# Quantitative Data

The performance of **PatMaN** has been benchmarked for tasks such as mapping microarray probes to a genome. The following table summarizes key performance metrics reported in the original publication.

Parameter	Value
CPU	2.2 GHz Workstation
RAM Usage	~260 MB
Task	Matching 201,807 Affymetrix HGU95-A 25mer probes to the chimpanzee genome (panTro2)
Allowed Mismatches	1
Allowed Gaps	0
Execution Time	~2.5 hours
Total Hits Found	15.9 million

## Experimental Protocols

**PatMaN** is a command-line tool, and its execution is controlled by a set of parameters. A typical experimental workflow involves preparing the input files, running the **PatMaN** executable with the desired options, and then processing the output.

### 4.1. Input Data Preparation

- Query Sequences: Create a FASTA file containing the short nucleotide sequences to be searched. Each sequence should have a unique identifier.
- Database Sequences: Prepare a FASTA file with the large nucleotide database (e.g., a chromosome or an entire genome).

### 4.2. Execution via Command Line

The basic command structure for running **PatMaN** is as follows:

Key Command-Line Options:

Option	Description
-P, --patterns	Specifies the input file containing the pattern (query) sequences in FASTA format.
-D, --databases	Specifies the input file containing the database sequences in FASTA format.
-e, --edits	Sets the maximum number of edits (mismatches + gaps) allowed per match.
-g, --gaps	Sets the maximum number of gaps allowed per match. Note that gaps also count as edits.
-o, --output	Redirects the output to the specified file. The default is standard output.
-a, --ambicodes	Activates the interpretation of ambiguity codes in the pattern sequences.
-s, --singlestrand	Deactivates matching of the reverse-complements of the patterns.

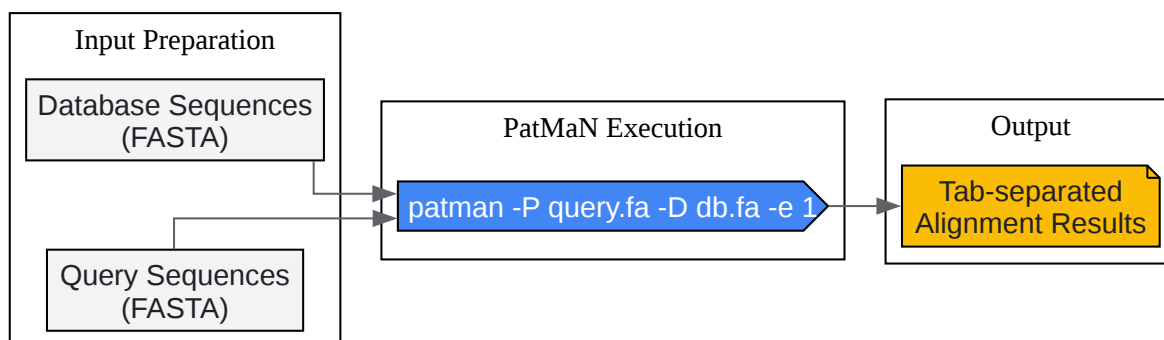
### 4.3. Output Format

**PatMaN** produces a tab-separated output file with the following columns for each match found:

- Database sequence name
- Pattern name
- Start position of the match in the database sequence (1-based)
- End position of the match in the database sequence
- Strand (+ for forward, - for reverse complement)
- Edit distance (number of mismatches + gaps)

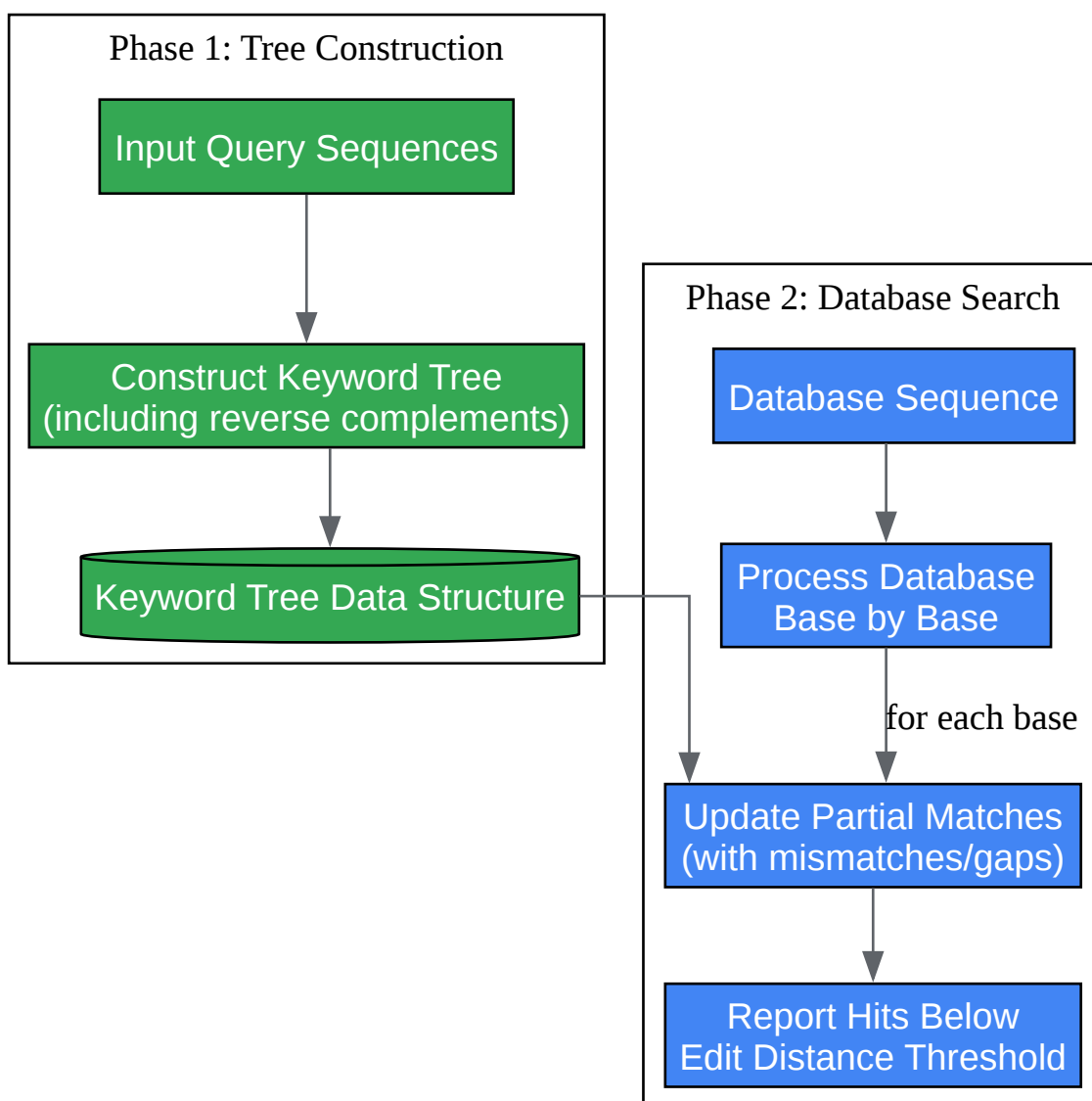
## Visualizations

To further elucidate the functionality of **PatMaN**, the following diagrams illustrate the experimental workflow and the core algorithmic logic.



[Click to download full resolution via product page](#)

A high-level overview of the **PatMaN** experimental workflow.



[Click to download full resolution via product page](#)

The logical flow of the core **PatMaN** algorithm.

#### Need Custom Synthesis?

BenchChem offers custom synthesis for rare earth carbides and specific isotopic labeling.

Email: [info@benchchem.com](mailto:info@benchchem.com) or [Request Quote Online](#).

## References

- 1. PatMaN: rapid alignment of short sequences to large databases - PMC [pmc.ncbi.nlm.nih.gov]
- 2. academic.oup.com [academic.oup.com]
- 3. PatMaN - Bioinformatics DB [bioinformaticshome.com]
- To cite this document: BenchChem. [PatMaN: A Technical Guide to a High-Throughput Short Sequence Search Tool]. BenchChem, [2025]. [Online PDF]. Available at: [https://www.benchchem.com/product/b1221989#patman-source-code-and-documentation]

---

### Disclaimer & Data Validity:

The information provided in this document is for Research Use Only (RUO) and is strictly not intended for diagnostic or therapeutic procedures. While BenchChem strives to provide accurate protocols, we make no warranties, express or implied, regarding the fitness of this product for every specific experimental setup.

**Technical Support:** The protocols provided are for reference purposes. Unsure if this reagent suits your experiment? [[Contact our Ph.D. Support Team for a compatibility check](#)]

**Need Industrial/Bulk Grade?** [Request Custom Synthesis Quote](#)

## BenchChem

Our mission is to be the trusted global source of essential and advanced chemicals, empowering scientists and researchers to drive progress in science and industry.

### Contact

Address: 3281 E Guasti Rd

Ontario, CA 91761, United States

Phone: (601) 213-4426

Email: [info@benchchem.com](mailto:info@benchchem.com)