

## PAPER

# PxBLAT: An Ergonomic and Efficient Python Binding Library for BLAT

Yangyang Li<sup>1</sup> and Rendong Yang<sup>1,2,\*</sup><sup>1</sup>Department of Urology, Northwestern University, Feinberg School of Medicine, Chicago, IL, 60611, USA and <sup>2</sup>Robert H. Lurie Comprehensive Cancer Center, Northwestern University Feinberg School of Medicine, Chicago, IL, 60611, USA

\*Corresponding author. rendong.yang@northwestern.edu

FOR PUBLISHER ONLY Received on Date Month Year; revised on Date Month Year; accepted on Date Month Year

## Abstract

**Summary:** We introduce PxBLAT, a Python library designed to enhance usability and efficiency in interacting with the BLAT. PxBLAT provides an intuitive API design, allowing the incorporation of its functionality directly into Python-based bioinformatics workflows. Besides, it integrates seamlessly with Biopython and comes equipped with user-centric features like server readiness checks and port retry mechanisms. PxBLAT removes the necessity for system calls and intermediate files, as well as reducing latency and data conversion overhead. Benchmark tests show that PxBLAT outperforms BLAT in the Python environment by roughly 20%.

**Availability and implementation:** PxBLAT supports Python (version 3.8+), and pre-compiled packages are released via PyPI (<https://pypi.org/project/pxblat/>) and Bioconda (<https://anaconda.org/bioconda/pxblat>). The source code of PxBLAT is licensed under an open-source MIT license and is accessible via GitHub (<https://github.com/ylab-hi/pxblat>). Its documentation is available on ReadTheDocs (<https://pxblat.readthedocs.io/en/latest/>).

**Supplementary information:** Supplementary data are available at *Bioinformatics* online.

**Key words:** Software Libraries, Sequence Analysis, BLAT

## Introduction

The rise of Python as a preferred programming language within bioinformatics is widely acknowledged as a result of its user-friendly nature, extensive libraries, and unparalleled versatility (Perkel, 2015). A variety of libraries have been crafted to augment Python's interface, thereby amplifying the adaptability and compatibility of bioinformatics tools (Putri et al., 2022; Cock et al., 2009). For instance, Biopython (Cock et al., 2009), a preeminent bioinformatics library, furnishes interfaces to tools like BLAST (Altschul et al., 1990) and Clustal (Higgins and Sharp, 1988).

Additionally, the relentless progression of genome sequencing technologies has precipitated a dramatic surge in the availability of genomic data. The necessity for tools to decipher and manipulate these data has become paramount within both the research and clinical domains. BLAT (Kent, 2002) is a standout within the bioinformatics landscape and is recognized for its capability to swiftly conduct genome sequence alignments. It offers a faster alternative to BLAST (Altschul et al., 1990) for aligning DNA sequences to the human genome (Kent, 2002). Despite its widespread use and endorsement within the bioinformatics community, interfacing with BLAT can present challenges. BLAT, implemented in the C language, solely extends utilities via CLIs, rendering its integration into Python-based projects inconvenient. Moreover,

BLAT triggers system calls and data conversions that cause significant overhead. Until now, a comprehensive approach to enhancing the usability of BLAT has been absent when incorporated into Python-based projects.

Here, we propose PxBLAT, which empowers users to utilize BLAT programmatically, facilitating its seamless integration into novel algorithms or analytical pipelines. PxBLAT serves as a bridge, bringing the high-performance capabilities of BLAT into the Python environment while maintaining reproducibility. The overarching aim of this work is to fill the observed gap by providing a Python binding library specifically tailored for BLAT, addressing both its efficiency and ergonomic concerns.

## Implementation

The design of PxBLAT adheres to the principles of readability and simplicity, aiming for an intuitive user experience that diminishes the learning curve. Striving to curtail complexity while enhancing usability and performance, we extracted the implementation of BLAT from the codebase of the UCSC, thus reducing dependencies. We retain the original C codebase and reimplement the utilities of BLAT (V37.1) including *faTwoBit*, *gfServer*, and *gfClient* CLI, using the C++ programming language (Kent, 2002). Subsequently, we integrate the current C++ code and develop PxBLAT via Pybind11 (Jakob et al.,

**Table 1.** Performance Benchmarking

Group	Samples (No.)	BLAT (Sec.)	PxBLAT (Sec.)	Speedup
1	200.00	61.40	50.41	1.22
2	200.00	61.70	51.10	1.21
3	200.00	55.34	46.44	1.19
4	200.00	56.98	47.72	1.19
5	200.00	54.55	46.44	1.17

2016). The approach facilitated direct interaction with the functions of the BLAT without modifying its original source code, thereby preserving the integrity and performance of the BLAT while broadening its capabilities.

The query result of PxBLAT conforms to the *QueryResult* class of Biopython (Cock et al., 2009), enabling manipulation of the query result using Biopython features (Listing 1). Furthermore, PxBLAT eliminates the need for intermediate files, allowing all operations to be executed in memory. This eradicates the common obstacle of data conversion into specific formats, enabling users to focus more on the core sequence alignment task. Input and output files have been made optional, offering more flexible choices. The use of system calls, though functional, can induce latency and create performance bottlenecks. PxBLAT reduces the need for system calls, thereby improving efficiency (Table 1). Moreover, PxBLAT allows for the retrieval of server status without manipulating log files, a process that can be troublesome in concurrent environments. The library incorporates ergonomic features such as checking and waiting for server readiness for alignment, retrying different ports if the current one is occupied, and using an existing server if one has already been started.

We provide a variety of examples and documentation to help the user get started (Listing 1). Specifically, PxBLAT provides APIs, including Class *Server*, *Client*, and other free functions, to reproduce the BLAT suit. Class *Server* and *Client* have the same utilities as CLI *gfServer* and *gfClient*, respectively, but with greater flexibility. Free functions, for example, *start\_server*, *query\_server*, *status\_server*, *fa2twobit*, *twobit2fa*, etc., serve as potential usage in different contexts. The library has been tested and developed with CI and CD to ensure code quality. It uses type annotations for public classes and functions to guarantee quality and correctness via a type checker and static analyzer. The type annotations also make PxBLAT more user-friendly in developing environments, where the function signatures can be suggested and corrected automatically. Besides, PxBLAT contains CLIs implemented by its API. The CLI includes completion for different shells, improving versatility.

## Benchmarking

The result and performance of PxBLAT are benchmarked against BLAT (V37.1) using a dataset of 1000 FASTA Files. The testing dataset, which consists of five groups, is randomly sampled from chromosome 20 of the Human Genome (hg38), and each sample includes one sequence. The length of each sequence ranges from 1000 bp to 3000 bp, which covers a range of typical usage scenarios (Figure S1). We summarized the HSPs of all hits in each sample derived from BLAT and PxBLAT and compared their predictions. We observed that PxBLAT reproduces the HSPs of BLAT across all tested samples, which ensures the correctness of PxBLAT (Tables S1 to S5).

**Listing 1** API Example The code snippet shows how to use the API of PxBLAT, and the query result can be iterated. More code examples can be found at <https://pxblat.readthedocs.io/en>

```

from pxblat import Server
from pxblat import Client

client = Client(
    host="localhost",
    port=65000,
    seq_dir="ref/",
    min_score=20,
    min_identity=90,
)

server_option = Server.create_option().build()
with Server(
    host="localhost",
    port=65000,
    two_bit="ref/reference.2bit",
    option=server_option
) as server:
    work() # work that consumes time
    server.wait_for_ready()
    result1 = client.query("ATCG")
    result2 = client.query("AtcG")
    result3 = client.query(["ATCG", "ATCG"])
    result4 = client.query(["cgTA", "fasta.fa"])

for hsp in result1.hsps:
    print(f"{hsp}")

```

The benchmarking was conducted on an Apple M1 Pro with macOS 13.4.1 22F82 arm64. We use system calls to launch BLAT and the *time* library to measure the execution time. PxBLAT, on average, gains 20% speedup compared to BLAT (Table 1). In summary, PxBLAT offers tangible benefits according to reduced execution time and improved user experience, proving its value as an enhancement to BLAT's functionality.

## Acknowledgments

Special thanks to the team who maintain the UCSC codebase and users from the bioinformatics community whose valuable feedback and suggestions were pivotal in refining PxBLAT's design and functionality.

## Competing interests

No competing interest is declared.

## Funding

This work was supported by the National Institute of General Medical Sciences [R35GM142441].

## Data availability

The PxBLAT, along with the source code, is publicly available in the GitHub repository at <https://github.com/yliab-hi/>

pxblat. The documentation is available at ReadtheDocs <https://pxblat.readthedocs.io/en/latest/>. The script for benchmarking is available at `tests/test_result.py` in the repository. The testing dataset is available at the GitHub repository <https://github.com/ylab-hi/pxblat>. The path of the testing dataset is `benchmark/fas`.

## References

- S. F. Altschul, W. Gish, W. Miller, E. W. Myers, and D. J. Lipman. Basic local alignment search tool. *Journal of Molecular Biology*, 215(3):403–410, 1990. ISSN 0022-2836. doi: 10.1016/S0022-2836(05)80360-2.
- P. J. A. Cock, T. Antao, J. T. Chang, B. A. Chapman, C. J. Cox, A. Dalke, I. Friedberg, T. Hamelryck, F. Kauff, B. Wilczynski, and M. J. L. de Hoon. Biopython: freely available python tools for computational molecular biology and bioinformatics. *Bioinformatics*, 25(11):1422–1423, 2009. doi: 10.1093/bioinformatics/btp163. URL <https://doi.org/10.1093/bioinformatics/btp163>.
- D. G. Higgins and P. M. Sharp. CLUSTAL: A package for performing multiple sequence alignment on a microcomputer. *Gene*, 73(1):237–244, 1988. ISSN 0378-1119. doi: 10.1016/0378-1119(88)90330-7.
- W. Jakob, J. Rhinelanders, and D. Moldovan. pybind11 — seamless operability between c++11 and python, 2016. <https://github.com/pybind/pybind11>.
- W. J. Kent. BLAT—The BLAST-Like Alignment Tool. *Genome Research*, 12(4):656–664, 2002. ISSN 1088-9051, 1549-5469. doi: 10.1101/gr.229202.
- J. M. Perkel. Programming: Pick up Python. *Nature*, 518(7537):125–126, 2015. ISSN 1476-4687. doi: 10.1038/518125a.
- G. H. Putri, S. Anders, P. T. Pyl, J. E. Pimanda, and F. Zanini. Analysing high-throughput sequencing data in Python with HTSeq 2.0. *Bioinformatics*, 38(10):2943–2945, 2022. ISSN 1367-4803. doi: 10.1093/bioinformatics/btac166.