

Runtime Code Generation & Multi-Stage Programming in Python

[VICTOR KAMEL]

Programmers often appreciate the convenience of writing code in high-level, interpreted languages such as Python. However, the high level of abstractions in these languages conspire against the programmer by making the resulting code slow. When performance is required, this leads to the so-called “Two-Language Problem,” where part of the code is written in a slow-yet-flexible language such as Python, and the performance-critical sections are written in a fast-yet-rigid language like C.

Fortunately, techniques such as multi-stage programming and just-in-time (JIT) compilation allow code written in interpreted languages to reclaim some of the lost performance by partially compiling some of the code at runtime. While many Python implementation and libraries aim to provide JIT compilation, there are very few that provide facilities for multi-stage programming. In this project, we investigate the performance benefits of implementing a prototype of runtime code generation and multi-stage programming into Python, demonstrating a mean speedup of up to 25× over normal Python functions. In addition, we demonstrate how multi-stage programming can be used to create a compiler for a simple tape-based programming language.

1 Introduction

Multi-stage programming (MSP) allows—with the help of special annotations—programming languages to specialize (compile) code at runtime in a type-safe manner. Thus, we can gain the advantages of compiled and optimized code in an interpreted language.¹ This also allows for programming techniques such as generative programming, where we can use staged compilation to craft programs such as compilers for domain-specific languages (DSLs) that can avoid part of the performance loss incurred by using the abundance of abstractions available in high-level languages.

In this project, we implement a prototype that enables runtime compilation and multi-stage programming in Python in order to explore its performance benefits. Since Python is an incredibly popular language, there are many projects that aim to accelerate Python code via JIT compilation, beyond a simple foreign function interface (FFI), explored in §5. We leverage similar methods to achieve runtime compilation. However, unlike those approaches, we implement staged compilation in order to enable generative programming. Our prototype² supports the following features, enabled via various mechanisms:

- (1) compilation and optimization of functions written in a subset of Python to native code at runtime;
- (2) multistage programming with splices and quotes;
- (3) integer, float, boolean, and multi-dimensional array datatypes;
- (4) basic control flow, labels, and goto statements; and
- (5) calling external C functions;

Although Python is dynamically typed, since compilation now occurs at runtime, we can generate properly-typed code that is interoperable with the running Python instance as needed. We show that the generated functions offer a significant speedup over pure Python, and demonstrate that MSP can be used to create a compiler for a DSL.

¹If we can amortize the cost of compilation, and the compiled code comprises a significant proportion of the runtime (Amdahl’s law).

²Source code available: <https://github.com/HeroHFM/msp>

2 Methodology

In this section, we present the design of our Python MSP module, and details of the corresponding implementation.

2.1 Compilation

In order to compile a Python function at runtime, we provide the `compile` decorator. All function parameters as well as the return type must have a type annotation. In order to compile the function, the following must occur:

- (1) the source code of the function is retrieved using the `inspect`³ module's `getsource` function;
- (2) the source code is converted to an abstract syntax tree (AST) representation using the `ast`⁴ module's `parse` function;
- (3) the AST is translated into LLVM IR using a custom `ast.NodeVisitor` and Numba's `llvmlite`⁵ binding layer;
- (4) the AST is used to determine the signature of the function; and
- (5) the compiled function along with its signature are saved in memory.

Subsequently, `ctypes`⁶ can be used to run the function from Python using Python datatypes as input arguments. We chose to use LLVM intermediate representation (IR), this allows us to leverage LLVM's optimization passes to optimize the generated code automatically. The following is an example of compiling and running a simple function using our library:

```

1 from msp.msp import *
2 gen = CGEngine.create_engine()
3
4 @gen.compile()
5 def add3(n : int) -> int:
6     return n + 1 + 1 + 1
7
8 # >>> add3(0)
9 # 3

```

The first 2 lines import our module, and create a code generation engine object that will be subsequently used to compile code. The `compile` decorator automatically performs all steps outlined above, and makes the function available to call in the current environment. The AST of `add3` and the unoptimized LLVM IR translation is available in Appendix A. The optimized LLVM code is as follows:

```

1 ; Function Attrs: mustprogress norecurse nosync nounwind readnone willreturn
2 define i32 @add3(i32 %n) local_unnamed_addr #0 {
3 entry:
4     %.5 = add i32 %n, 3
5     ret i32 %.5
6 }
7
8 attributes #0 = { mustprogress norecurse nosync nounwind readnone willreturn }

```

Note that the three additions were collapsed into a single operation. In general, LLVM can perform much more sophisticated optimizations than this given more complex code.

³<https://docs.python.org/3/library/inspect.html>

⁴<https://docs.python.org/3/library/ast.html>

⁵<https://github.com/numba/llvmlite>

⁶<https://docs.python.org/3/library/ctypes.html>

2.2 Multi-stage programming

MSP⁷ typically uses the following code annotations in order to enable staging: quotes, splices, and a mechanism for running code. Quotes allow us to represent code fragments as a unit that can be manipulated (e.g. returned from a function), splices allow us to “splice in” quoted code into our function, and finally we require a mechanism to compile and run the code.

2.2.1 Quote. To enable quotation, we provide the user with the function `q`, which takes the code in the form of a string as a parameter. We chose to make the code inside of the quotation a string, as it allows more flexible manipulation of the quotation using Python. The function evaluates to a Python AST node via the `eval` function.

2.2.2 Splice. To enable splices, we define the `splice` function⁸. This function takes a regular Python expression, and evaluates it using Python’s `eval` keyword. If this produces an AST node, the node gets spliced into the AST of the calling function. If the splice operation contains a function call, we must also perform a rewrite pass in order to replace the parameters with the appropriate arguments it was called with.

2.2.3 Run. Any function that is defined using the `compile` decorator is automatically compiled. When this function is subsequently called, the compiled version of the function is executed. The following example demonstrates these mechanisms in action.

```

1 def power(n, x):
2     if n == 0: return q('1')
3     return q('splice(power(n - 1, x)) * splice(x)')
4
5 @gen.compile(env = globals())
6 def power10(x : int) -> int:
7     return splice(power(10, q('x')))
8
9 # >>> power10(2)
10 # 1024

```

This code specializes the `power` function at runtime. Note that the `power` function is a regular Python function that returns an AST node (since both returns contain a quote). Since this function is not compiled and runs only at runtime, there is no restriction on what code it may contain. In this instance, the `compile` decorator must take an additional argument—the environment—in order to “capture” the definition of the `power` function. The return value splices in the AST node retrieved from `power`, where the first argument is the value 10, and the second is the quotation of `x`. Once the first node is returned, the resulting expression is evaluated, and the subsequent splices are also evaluated. Finally, we note that the quotation of `x` is used in the call to `power`, since the value of `x` is not known at compile time, and thus we cannot evaluate `power` with `x` directly. Now that we have specialized the function, we can generate optimized code:

```

1 ; Function Attrs: mustprogress nofree norecurse nosync nounwind readonly willreturn
2 define i32 @power10(i32 %x) local_unnamed_addr #0 {
3     entry:
4         %0 = mul i32 %x, %x
5         %1 = mul i32 %0, %0
6         %2 = mul i32 %1, %x
7         %3 = mul i32 %2, %2
8         ret i32 %3

```

⁷For a more complete introduction to multi-stage programming, see [Taha 2004].

⁸Available within our compiled functions only.

```

9 }
10
11 attributes #0 = { mustprogress nofree norecurse nosync nounwind readnone willreturn }

```

2.3 Control Flow

We support several types of control flow, including If statements and While statements. More Python control flow structures can be added without much difficulty, but these were selected due to their versatility. In addition, `goto` and `label` were added as calls available in compiled functions, in order to aid with the implementation of the compiler in §3.2. In the following function definition, the `goto` call is used to break out of the While loop.

```

1 @gen.compile()
2 def double_positive_numbers(n : int) -> int:
3     i = 0
4     acc = 0
5     while True:
6         if i >= n: goto("exit")
7         else: i = i + 1
8         acc = acc + 2
9     label("exit")
10    return acc
11
12 # >>> double_positive_numbers(5)
13 # 10

```

2.4 Datatypes

Our implementation mostly has support for numeric datatypes, including integers, floats, booleans. This is due to the fact that `llvmlite` was developed for Numba⁹, which is a JIT compiler for numerical computations. In addition, we have support for multi-dimensional arrays, as long as they are of a single type and their extents are known at compile time. In order to call these functions, the arrays must be created with Numpy¹⁰. The Numpy arrays are automatically marshalled by our backend (cast to a pointer type). We support loading from and storing to multi-dimensional arrays.

```

1 import numpy as np
2
3 array = np.array([[1, 2, 3], [4, 5, 6], [7, 8, 9]], dtype=np.int32)
4
5 @gen.compile(env = globals())
6 def sum_ndarray(arr : Int[3][3]) -> int:
7     acc = 0
8     i = 0
9     while i < 3:
10        j = 0
11        while j < 3:
12            acc = acc + arr[i, j]
13            j = j + 1
14        i = i + 1
15    return acc
16
17 # >>> sum_ndarray(array)
18 # 45

```

⁹<https://numba.pydata.org/>

¹⁰<https://numpy.org/>

2.5 Calling external functions

We permit calling external functions (such as those from `libc`), provided that a stub of their function signature is provided (passed to the `compile` decorator). For example, here we use the `putchar` function from the C standard library to print out a matrix of characters.

```

1 array = np.array([[72, 105, 32], [116, 104, 101], [114, 101, 33]], dtype=np.int32)
2
3 def putchar(c : int) -> int: pass
4
5 @gen.compile(env = globals(), external = {putchar,})
6 def sum_ndarray(arr : Int[3][3]) -> None:
7     acc = 0
8     i = 0
9     while i < 3:
10         j = 0
11         while j < 3:
12             putchar(arr[i, j])
13             putchar(32) # " "
14             j = j + 1
15             putchar(10) # "\n"
16             i = i + 1
17         return
18
19 # >>> sum_ndarray(array)
20 # H i
21 # t h e
22 # r e !

```

2.6 Scope extrusion

It may be of some interest whether our implementation permits scope extrusion, such as that described in [Calcagno et al. 2003]. We can implement a similar example:

```

1 l = 1
2
3 def inner():
4     global l
5     l = q('x + 1')
6     return q('2')
7
8 @gen.compile(env = globals())
9 def f(x : int) -> int:
10     return splice(inner())
11
12 # >>> l
13 # <ast.Module at 0x7f78aa217a90>

```

Indeed, the quote has escaped¹¹. However, we are not able to run `l` directly (unless we manually use Python's `compile` and `eval`), since we can only run compiled functions. We do not permit global functions in compiled functions either, so placing it in a function would fail to compile. However, this will not lead to a crash. At most, we will get a runtime exception (since compilation occurs at runtime), which can be handled.

¹¹We could just as easily have simply assigned the quote to `l` directly to achieve the same effect. Python is dynamically typed, so we cannot rely on the type system here.

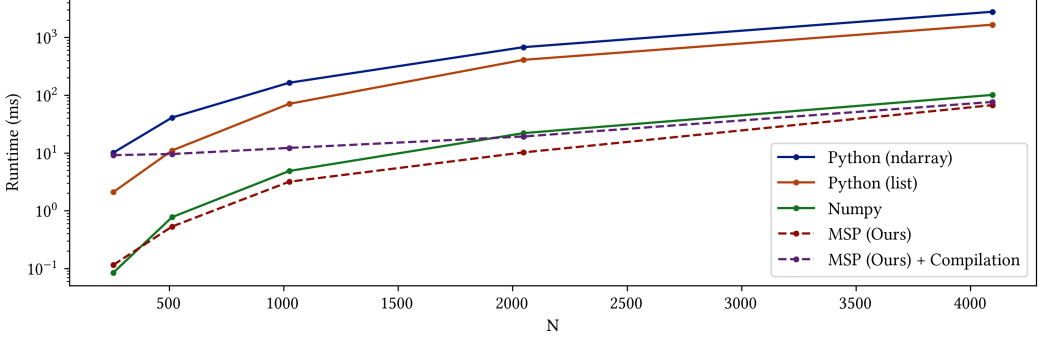


Fig. 1. Transpose running time on $N \times N$, $N \in \{2^8, 2^9, 2^{10}, 2^{11}, 2^{12}\}$ matrix.

3 Results

In this section, we will present some performance comparisons between normal Python functions and compiled functions, as well as to showcase an example of using MSP to compile a simple DSL, demonstrating the added benefit of MSP over plain JIT compilation.

3.1 Performance

It is clear that a compiled function will run faster than an interpreted function. However, since there is compilation overhead, a more detailed look at the performance breakdown is required. We will use a simple matrix transposition function to determine the speed advantage. Performance numbers were collected with the IPython `%timeit` line magic¹².

3.1.1 Matrix Transposition. We perform a transposition of a $N \times N$ matrix of floating point numbers. The functions used are available in Appendix B. Two pure-Python functions are compared. One uses the default `list` datatype, while the other uses a Numpy array. We observe that the naive use of the Numpy array does not improve the performance of the Python function. In addition, we test the use of a Numpy function to transpose the array¹³, which achieves a significant speedup. Our implementation achieves an average speedup of $25\times$ over the fastest Python version, and $1.5\times$ over the Numpy implementation. However, compiling the MSP function takes an average of 9.03ms. Therefore, if the function call is used only once (or the function runs for too short a time), Numpy may be the better choice until N becomes sufficiently large.

3.2 DSL Compilation

§3.1 demonstrated that there can be a significant benefit to compiling Python functions to native code. However, we want to further demonstrate the usefulness of the MSP approach. To do so, we write a compiler for the `bf`¹⁴ language, which is a Turing complete programming language resembling the operation of a Turing machine. In `bf`, we have a tape comprised of integer cells and a pointer. All cells initially hold the value zero, and the pointer begins at the leftmost cell. The `>` and `<` commands move the pointer right and left along the tape respectively, while the `+` and `-` commands increment or decrement the value of the current cell. The `.` command prints the value of the current cell in ASCII, and `,` prompts the user for a character and writes the corresponding

¹²<https://ipython.readthedocs.io/en/stable/interactive/magics.html>

¹³By default, Numpy's transposition creates a view and does not move any data. Therefore, we copy the array to ensure fairness between implementations.

¹⁴<https://esolangs.org/wiki/Brainfuck>

ASCII value to the cell. The `[` and `]` commands can be used to create loops. We can construct our compiler as follows, with source code that prints a greeting:

```

1  from msp.msp import *
2  import numpy as np
3
4  gen = CGEngine.create_engine()
5
6  source = "+++++++[>+++++>" \
7           "+++++++>+><<<" \
8           "-]>+.>+.+++++.++" \
9           "+.>+.<<+++++++++" \
10          "++>+.++-.-.-.-" \
11          "-.-.-.>+.>."
12
13 def putchar(n : int) -> int: pass
14 def getchar() -> int: pass
15
16 def bf():
17     program = []
18     jmp = 0
19     stack = []
20     for command in source:
21         if command == '+': program.append('data[ptr] = data[ptr] + 1')
22         elif command == '-': program.append('data[ptr] = data[ptr] - 1')
23         elif command == '>': program.append('ptr = ptr + 1')
24         elif command == '<': program.append('ptr = ptr - 1')
25         elif command == '.': program.append('_ = putchar(data[ptr])')
26         elif command == ',': program.append('data[ptr] = getchar()')
27         elif command == '[':
28             before, after = f"pre_{jmp}", f"post_{jmp}"
29             stack.append((before, after))
30             jmp += 1
31             program.append("\n".join([
32                 f"label('{before}')" ,
33                 f"if data[ptr] == 0: goto('{after}')" ,
34             ]))
35         elif command == ']':
36             before, after = stack.pop()
37             program.append("\n".join([
38                 f"goto('{before}')" ,
39                 f"label('{after}')" ,
40             ]))
41     return q("\n".join(program))
42
43     return body
44
45 times = []
46 @gen.compile(env = globals(), debug = False, external = {putchar, getchar})
47 def run(data : Int[256]) -> None:
48     ptr = 0
49     splice(bf())
50     return
51
52 arr = np.zeros(30000, dtype=np.int32)
53
54 run(arr)

```

We can see that the `run` function simply splices in its body from the `bf` function, which generates a quote based on the `bf` source code. For each command, the corresponding code is added into the program. Where jumps are required, the `label` and `goto` statements are used. This allows the jumps to be pre-calculated in the resulting IR. While this produces quite long source code, LLVM can optimize it to a certain extent. Upon compilation, the `run` function contains the code required to print “Hello World!”

4 Discussion

The current implementation is still at a prototype stage, with many parts of Python syntax still left unimplemented. Adding additional control flow, datatypes, and other Python features, such as generators and list comprehensions, can be considered future work. However, since there are already several projects that aim to JIT-compile Python, a better approach may be to extend an existing JIT compiler with staging annotations. In addition, many of these implementations focus on numeric data (including our own), so it would be interesting to implement other datatypes such as strings, which would allow replicating [Rompf and Amin \[2015\]](#).

5 Related work

There are several projects that attempt compilation in Python, some of which integrate with CPython (the reference implementation of Python), and some that are substitutes for it. We organize them into several categories, and only describe a representative of each category for brevity.

Python compilers. Cython¹⁵ is a Python-to-C compiler that can interoperate with CPython.

CPython alternatives. PyPy¹⁶ is a faster implementation of Python that includes a tracing JIT-compiler.

Numerical accelerators. Libraries such as Numba¹⁷ can compile and JIT-compile numerical Python at runtime.

JIT compilers. Pyjion¹⁸ is a JIT extension for CPython that compiles to Microsoft’s CIL.

DSLs. PyTorch¹⁹, Triton²⁰ and others implement DSLs that can be compiled to optimized CPU/GPU code for machine learning. Other such DSLs exist for other fields as well, such as Seq²¹ for genomics (now Codon [\[Shajii et al. 2023\]](#)), a more general framework and compiler for high-performance DSLs).

Multi-stage programming. Scale²² is an implementation of multi-stage programming, but it is not actively developed and has not been updated in many years.

Our approach aligns with that of Numba, as we take a Python function at runtime and lower it to LLVM IR. This is similar to the strategy employed by Scale. However, our splicing mechanism is different than theirs (we evaluate the function in Python, while they built their interpreter to do so). In addition, our multidimensional arrays are implemented differently. Our implementation is consistent with Numpy arrays, while theirs is implemented like Python lists, requiring expensive copy operations to move data between Python lists and their compiled function.

There is a wide literature describing the use of multi-staged compilation for programming languages, typically with functional languages such as ML [\[Taha and Sheard 2000\]](#), Haskell [\[Sheard](#)

¹⁵<https://github.com/python/cpython>

¹⁶<https://pypy.org/>

¹⁷<https://numba.pydata.org/>

¹⁸<https://github.com/tonybaloney/Pyjion>

¹⁹<https://pytorch.org/>

²⁰<https://github.com/triton-lang/triton>

²¹<https://seq-lang.org/>

²²<https://github.com/es1024/python-staged-programming>

and Jones 2002] or Scala [Rompf and Odersky 2010]. In particular, Rompf and Amin [2015] claim that with their library-based Scala approach, an SQL compiler for in-memory data is competitive with a C implementation without needing to compromise on abstraction. They base their implementation on the idea of *Futamura Projections* [Futamura 1999], whereby a an interpreter that can specialize to a particular source code can be considered a compiler. This code specialization (in the form of AST manipulation followed by compilation) at runtime is exactly what we achieve in Python in this project.

6 Conclusion

In this project, we developed a prototype of a Python module that enables both runtime code generation through LLVM IR and multi-staged programming in Python. This allows us to effectively mitigate the the performance penalty of using Python compared to a more low-level language without giving up on any convenience, as well as enables easy creation of compilers for performant DSLs. We demonstrate the effectiveness of our approach by showing a 25× mean speedup over the pure Python version of a transpose function, as well as a 1.5× speedup over the Numpy version. In addition, we show that we can use the MSP functionality to create a compiler for a simple language.

References

- C. Calcagno, E. Moggi, and T. Sheard. 2003. Closed types for a safe imperative MetaML. *J. Funct. Program.* 13, 3 (May 2003), 545–571. <https://doi.org/10.1017/S0956796802004598>
- Yoshihiko Futamura. 1999. Partial Evaluation of Computation Process, Revisited. *Higher Order Symbol. Comput.* 12, 4 (Dec. 1999), 377–380. <https://doi.org/10.1023/A:1010043619517>
- Tiark Rompf and Nada Amin. 2015. Functional pearl: a SQL to C compiler in 500 lines of code. *SIGPLAN Not.* 50, 9 (Aug. 2015), 2–9. <https://doi.org/10.1145/2858949.2784760>
- Tiark Rompf and Martin Odersky. 2010. Lightweight modular staging: a pragmatic approach to runtime code generation and compiled DSLs. In *Proceedings of the Ninth International Conference on Generative Programming and Component Engineering* (Eindhoven, The Netherlands) (GPCE '10). Association for Computing Machinery, New York, NY, USA, 127–136. <https://doi.org/10.1145/1868294.1868314>
- Ariya Shajii, Gabriel Ramirez, Haris Smajlović, Jessica Ray, Bonnie Berger, Saman Amarasinghe, and Ibrahim Numanagić. 2023. Codon: A Compiler for High-Performance Pythonic Applications and DSLs. In *Proceedings of the 32nd ACM SIGPLAN International Conference on Compiler Construction* (Montréal, QC, Canada) (CC 2023). Association for Computing Machinery, New York, NY, USA, 191–202. <https://doi.org/10.1145/3578360.3580275>
- Tim Sheard and Simon Peyton Jones. 2002. Template meta-programming for Haskell. In *Proceedings of the 2002 ACM SIGPLAN Workshop on Haskell* (Pittsburgh, Pennsylvania) (Haskell '02). Association for Computing Machinery, New York, NY, USA, 1–16. <https://doi.org/10.1145/581690.581691>
- Walid Taha. 2004. *A Gentle Introduction to Multi-stage Programming*. Springer Berlin Heidelberg, Berlin, Heidelberg, 30–50. https://doi.org/10.1007/978-3-540-25935-0_3
- Walid Taha and Tim Sheard. 2000. MetaML and multi-stage programming with explicit annotations. *Theoretical Computer Science* 248, 1 (2000), 211–242. [https://doi.org/10.1016/S0304-3975\(00\)00053-0](https://doi.org/10.1016/S0304-3975(00)00053-0) PEPM'97.

A AST & LLVM Translation of `add3`

AST representation:

```

1 Module(
2   body=[
3     FunctionDef(
4       name='add3',
5       args=arguments(
6         posonlyargs=[],
7         args=[
8           arg(
9             arg='n',
10            annotation=Name(id='int', ctx=Load()))],
11       kwnonlyargs=[],
12       kw_defaults=[],
13       defaults=[]),
14     body=[
15       Return(
16         value=BinOp(
17           left=BinOp(
18             left=BinOp(
19               left=Name(id='n', ctx=Load()),
20               op=Add(),
21               right=Constant(value=1)),
22             op=Add(),
23             right=Constant(value=1)),
24           op=Add(),
25           right=Constant(value=1))),
26       decorator_list=[],
27       returns=Name(id='int', ctx=Load()),
28       type_params=[]],
29     type_ignores=[])

```

Unoptimized LLVM translation:

```

1 ; ModuleID = ""
2 target triple = "unknown-unknown-unknown"
3 target datalayout = ""
4
5 define i32 @"add3"(i32 %"n")
6 {
7   entry:
8     %".3" = add i32 %"n", 1
9     %".4" = add i32 %".3", 1
10    %".5" = add i32 %".4", 1
11    ret i32 %".5"
12 }

```

B Transpose Functions

Python function with np.ndarray datatype:

```

1 %%timeit A, B = np.random.rand(n, n).astype(np.float32), np.zeros((n, n), dtype=np.float32)
2
3 # Transpose Python (Numpy array)
4 def transpose_py1(array, out):
5     for i in range(n):
6         for j in range(i + 1, n):
7             out[i, j] = array[j, i]
8             out[j, i] = array[i, j]
9
10 transpose_py1(A, B)

```

Python function with list datatype:

```

1 %%timeit A, B = np.random.rand(n, n).astype(np.float32).tolist(), np.zeros((n, n), dtype=np.
    ↳ float32).tolist()
2
3 # Transpose Python (Python list)
4 def transpose_py2(array, out):
5     for i in range(n):
6         for j in range(i + 1, n):
7             out[i][j] = array[j][i]
8             out[j][i] = array[i][j]
9
10 transpose_py2(A, B)

```

Numpy function:

```

1 %%timeit A, B = np.random.rand(n, n).astype(np.float32), np.zeros((n, n), dtype=np.float32)
2
3 B = A.T.copy()

```

MSP function:

```

1 @gen.compile(env = globals(), debug = False, add_time = times)
2 def transpose_msp(array : Float[n][n], out : Float[n][n], n : int) -> None:
3     i = 0
4     while i < n:
5         j = i + 1
6         while j < n:
7             out[i, j] = array[j, i]
8             out[j, i] = array[i, j]
9             j = j + 1
10        i = i + 1
11    return

```

```

1 %%timeit A, B = np.random.rand(n, n).astype(np.float32), np.zeros((n, n), dtype=np.float32)
2
3 transpose_msp(A, B, n)

```