

# 1 Polyglot Jet Finding

2 *Graeme Andrew Stewart*<sup>1,\*</sup>, *Philippe Gras*<sup>2</sup>, *Benedikt Hegner*<sup>1</sup>, and *Atell Krasnopolski*<sup>3</sup>

3 <sup>1</sup>CERN, Esplanade des Particules 1, Geneva, Switzerland

4 <sup>2</sup>IRFU, CEA, Université Paris-Saclay, Gif-sur-Yvette, France

5 <sup>3</sup>Taras Shevchenko National University of Kyiv, Ukraine

6 **Abstract.** The evaluation of new computing languages for a large community,  
7 like HEP, involves comparison of many aspects of the languages' behaviour,  
8 ecosystem and interactions with other languages. In this paper we compare  
9 a number of languages using a common, yet non-trivial, HEP algorithm: the  
10 anti- $k_T$  clustering algorithm used for jet finding. We compare specifically the  
11 algorithm implemented in Python (pure Python and accelerated with numpy and  
12 numba), and Julia, with respect to the reference implementation in C++, from  
13 Fastjet. As well as the speed of the implementation we describe the ergonomics  
14 of the language for the coder, as well as the efforts required to achieve the best  
15 performance, which can directly impact on code readability and sustainability.

## 16 1 Introduction

17 High energy physics (HEP), as a discipline, has undergone at least two major shifts in lan-  
18 guage after the widespread adoption of Fortran in the 1960s [1]. A first was a significant shift  
19 from Fortran to C++, starting with the BaBar experiment, then gathering pace at the end of  
20 the the Large Electron-Positron Collider (LEP) era, c. 2000, when the Large Hadron Collider  
21 (LHC) experiments adopted C++ more or less wholesale. The second shift happened with  
22 the gradual incorporation of Python into the language ecosystem of HEP, from about 2010.

23 In the first transition, Fortran was almost completely displaced by C++ in the HEP experi-  
24 ments; in the theory domain the evolution was more gradual and mixed, with Fortran and C++  
25 still both used today. In the second, a different type of transition took place, where Python  
26 became more and more popular, but co-exists with C++. The C++ is largely used in per-  
27 formance critical areas, with Python finding traction when flexibility and rapid turn-around  
28 is needed, e.g., in configuration and steering. Python code is typically used to interface to  
29 higher performance C and C++ libraries, both generic (e.g., numpy) and specific HEP codes.

30 Although the field is, for reasons of stability and legacy, slow to move to new languages,  
31 there are some significant issues with the current language choices that make an exploration  
32 of alternatives worthwhile. For example, the interfaces between Python and C++ are a source  
33 of friction, both for passing data and error messages back and forth, as well as being obliged  
34 to switch languages and reimplement code on occasion, when moving from a prototype to  
35 production (assuming the the developer actually has skills in both languages, which is not  
36 a given). This *two language problem* has potentially been addressed in the *Julia* program-  
37 ming language [2, 3], with promising prospects for HEP, in particular, [4, 5] as well as other

---

\*e-mail: graeme.andrew.stewart@cern.ch

Language	Repository
C++ (FastJet)	FastJet Website (release 3.4.1)
Python (Pure)	GitHub antikt-python
Python (Accelerated)	GitHub antikt-python
Julia	GitHub JetReconstruction.jl

**Table 1.** Code repositories used in this paper. See [9] for exact commits and instructions.

38 STEM<sup>1</sup> areas [6]. Julia offers just in time compilation giving an ergonomic experience much  
39 like Python, but with runtime speeds comparable to C and C++. C++ is also a notoriously  
40 tricky language to use, particularly related to memory handling [7] and HEP C++ codes are  
41 frequently riddled with code defects [8].

42 Evaluation of the prospects for a language in any particular domain area should be done  
43 with a real problem from that domain, rather than any synthetic benchmark. In this paper we  
44 look at the problem of jet finding, or clustering, which is a use case from high energy physics  
45 used in calorimeter reconstruction. This is a good example as it is not trivial, but it is also not  
46 so complex that different implementations take too long to write.

47 The languages we examine here, along with links to the code used, are given in Table 1.

48 The evaluation itself can cover many aspects of a programming language and the expe-  
49 rience of using it. Metrics such as runtime are easy to evaluate, but the ergonomics of using  
50 particular languages and the support offered by the language ecosystem for developers are  
51 also critical and we comment on these.

## 52 2 Anti- $k_T$ Jet Clustering Algorithm

### 53 2.1 Algorithm

54 The anti- $k_T$  clustering algorithm [10, 11] is an infrared and colinear safe jet clustering algo-  
55 rithm, which is robust against soft fragmentation components. We use the Fastjet implemen-  
56 tation [12], that proceeds in the following way:

- 57 1. A radius parameter,  $R$ , is defined (0.4 is typical at the LHC).
- 58 2. For each active pseudojet  $A$  (that is, an initial particle or a merged cluster):
  - 59 (a) Considering all other PseudoJets,  $B$ , which are closer in geometric distance than  
60  $R$ , measure the minimum geometric distance:

$$61 \quad d = \min \left( \sqrt{\Delta\eta_{AB}^2 + \Delta\phi_{AB}^2} \right),$$

62 where  $\Delta\eta_{AB}$  and  $\Delta\phi_{AB}$  are the rapidity and azimuthal angle differences between  $A$   
63 and  $B$ .

64 If there are no other pseudojets within  $R$ , then  $d = R$  for pseudojet  $A$ .

- 65 (b) Define the anti- $k_T$  distance,  $d_{ij}$ , as  $d_{ij} = d \min(k_{T,A}^{-2}, k_{T,B}^{-2})$  where  $k_{T,\{A,B\}}$  is the  
66 transverse momentum of the pseudojet  $\{A, B\}$ . If there is no neighbouring pseu-  
67 dojet,  $d_{ij} = d k_{T,A}^{-2}$ .

- 68 3. Choose the pseudojet with the lowest  $d_{ij}$ .

---

<sup>1</sup>Science, technology, engineering, and mathematics

- 69 (a) If this pseudojet has an active partner,  $B$ , merge these two pseudojets to a new  
70 pseudojet.  
71 (b) If not, this jet is finalised and removed from the active list.  
72 4. Repeat until no pseudojets remain active.

73 Note that the definition of  $d_{ij}$ , so-called anti- $k_T$ , with a negative power favours merging  
74 jets with a high transverse momentum first, which provides stability against soft radiation,  
75 hence its popularity. (Considering a general metric distance of  $k_T^{2p}$ ,  $p = -1$  is anti- $k_T$  merging,  
76  $p = 0$  is Cambridge/Aachen merging and  $p = 1$  is inclusive  $k_T$ . [11].)

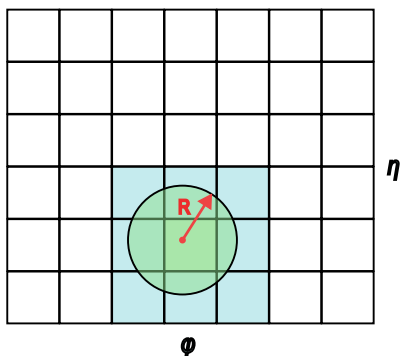
77 The algorithm itself has a nice mixture of parallelisation opportunities (pairwise matching  
78 of pseudojet candidates) and serial steps (finding the minimum values of  $d$  or  $d_{ij}$ ), which is a  
79 good test of a non-naive algorithm's performance.

## 80 2.2 Algorithm Implementations

81 We consider two different implementations of the algorithm described above, taken from  
82 FastJet [12, 13].

83 The first is a *plain implementation* in which, at each step, all jets are considered as possi-  
84 ble neighbours of each other. This algorithm has scaling that runs roughly as  $N^2$ , where  $N$  is  
85 the number of initial particle hits (this is an improvement over the most naive scaling which  
86 would be  $N^3$  [10]). This implementation is fastest for  $\lesssim 30$  particles.

87 The second is a *tiled implementation*, in which the geometric space  $(\eta, \phi)$  is split into tiles  
88 of size  $R$ . In this way the number of possible neighbours of any particular jet is limited to the  
89 jet's tile and to its immediate neighbours, as illustrated in Figure 1. This strategy reduces the  
90 amount of work that needs to be done, at the expense of extra bookkeeping of which jets are  
91 in each tile. This implementation is fastest for p-p collisions at the LHC.



**Figure 1.** In the tiled algorithm implementation  $(\eta, \phi)$  space is split into tiles of size  $R$ . When a pseudojet needs to rescan for neighbours (red dot) only pseudojets in tiles within the distance  $R$  need to be considered, here shaded in light blue.

## 92 3 Code Implementation Ergonomics

93 The code versions used are linked to in Table 1. We highlight some specific observations in  
94 this section.

## 95 3.1 C++, FastJet

96 The FastJet package [12, 13] is a well maintained code which is widely used in the HEP  
97 community. It is code which is of high quality and well scrutinised and tested. The general  
98 style of the code is more akin to C than C++, for reasons of minimising abstraction and  
99 increasing speed, although templates are used extensively (where any errors are not usually  
100 nicely handled by the compiler).

101 For the tiled implementation, a linked list structure is used, which requires pointers to  
102 pointers that are challenging to reason about for the programmer, as illustrated below.

```
// set up the initial nearest neighbour information
vector<Tile>::const_iterator tile;
for (tile = _tiles.begin(); tile != _tiles.end(); tile++) {
    // first do it on this tile
    for (jetA = tile->head; jetA != NULL; jetA = jetA->next) {
        for (jetB = tile->head; jetB != jetA; jetB = jetB->next) {
            double dist = _bj_dist(jetA,jetB);
            if (dist < jetA->NN_dist) {jetA->NN_dist = dist;
                jetA->NN = jetB;}
            if (dist < jetB->NN_dist) {jetB->NN_dist = dist;
                jetB->NN = jetA;}
        }
    }
    // then do it for RH tiles
    for (Tile ** RTile = tile->RH_tiles; RTile != tile->end_tiles; RTile++) {
        for (jetA = tile->head; jetA != NULL; jetA = jetA->next) {
            for (jetB = (*RTile)->head; jetB != NULL; jetB = jetB->next) {
                double dist = _bj_dist(jetA,jetB);
                if (dist < jetA->NN_dist) {jetA->NN_dist = dist;
                    jetA->NN = jetB;}
                if (dist < jetB->NN_dist) {jetB->NN_dist = dist;
                    jetB->NN = jetA;}
            }
        }
    }
}
```

## 103 3.2 Python

### 104 3.2.1 Pure Python

105 Python is renowned for being a high productivity language and implementation of the jet  
106 finding algorithms is rather straightforward, with a clear logic. Mutability of classes allows  
107 code to be shared between the different implementations. e.g., an update scan for the basic  
108 case looks like this:

```
def scan_for_my_nearest_neighbours(jetA: PseudoJet, jets: list[PseudoJet],
    R2: float):
    """Retest all other jets against the target jet"""
    jetA.info.nn = None
    jetA.info.nn_dist = R2
```

```

for ijetB, jetB in enumerate(jets):
    if not jetB.info.active:
        continue
    if ijetB == jetA.info.id:
        continue
    dist = geometric_distance(jetA, jetB)
    if dist < jetA.info.nn_dist:
        jetA.info.nn_dist = dist
        jetA.info.nn = ijetB
jetA.info.akt_dist = antikt_distance(jetA, jets[jetA.info.nn]
    if jetA.info.nn else None)

```

109 Where specifically `info` is a mix-in class for bookkeeping pseudojets.

110 While the code for the tiled implementation involves more bookkeeping, it also remains  
111 clear.

### 112 3.2.2 Accelerated Python

113 Accelerated Python code, where both numba and numpy are employed brings some added  
114 difficulty. Not all operations are easily expressed as numpy array calculations, particularly  
115 for dynamic arrays holding active and inactive jets. This necessitated the use of masks, which  
116 need to be tracked. In addition, numba jitted functions are very picky on types that can  
117 be passed (at least without being explicitly *taught* how to deal with them), so instead of a  
118 structure, functions are called with many individual array elements, leading to complicated  
119 call signatures. e.g., the same function as above becomes

```

@njit
def scan_for_my_nearest_neighbours(ijet:int, phi: npt.ArrayLike,
                                   rap:npt.ArrayLike, inv_pt2:npt.ArrayLike,
                                   dist:npt.ArrayLike, akt_dist:npt.ArrayLike,
                                   nn:npt.ArrayLike,
                                   mask:npt.ArrayLike, R2: float):
    "Retest all other jets against the target jet"
    nn[ijet] = -1
    dist[ijet] = R2
    _dphi = np.pi - np.abs(np.pi - np.abs(phi - phi[ijet]))
    _drap = rap - rap[ijet]
    _dist = _dphi*_dphi + _drap*_drap
    _dist[ijet] = R2 # Avoid measuring the distance 0 to myself!
    _dist[mask] = 1e20 # Don't consider any masked jets
    iclosejet = _dist.argmin()
    dist[ijet] = _dist[iclosejet]
    if iclosejet == ijet:
        nn[ijet] = -1
        akt_dist[ijet] = dist[ijet] * inv_pt2[ijet]
    else:
        nn[ijet] = iclosejet
        akt_dist[ijet] = dist[ijet] * (inv_pt2[ijet] if inv_pt2[ijet]
            < inv_pt2[iclosejet] else inv_pt2[iclosejet])
    # As this function is called on new PseudoJets it's possible

```

```

# that we are now the NN of our NN
if dist[iclosejet] > dist[ijet]:
    dist[iclosejet] = dist[ijet]
    nn[iclosejet] = ijet
    akt_dist[iclosejet] = dist[iclosejet] * (inv_pt2[ijet]
        if inv_pt2[ijet] < inv_pt2[iclosejet] else inv_pt2[iclosejet])

```

120 numba also has some surprising omissions from the numpy functions which it can JIT,  
 121 e.g., array index raveling, that required explicit reimplementaion.

### 122 3.3 Julia

123 Julia is gaining in popularity because it is a language that is easy to use. We found numerous  
 124 nice features that allow code to be clear, e.g., using the broadcast syntax for calculations on  
 125 arrays is very compact:

```
kt2 = (JetReconstruction.pt.(objects) .^ 2) .^ p
```

126 Here .^ (raise to power) operates on each member of the pt value of the objects array.  
 127 Like the FastJet code, loops can be used without sacrificing speed, so the code checking  
 128 for new nearest neighbours is

```

# Finds new nearest neighbour for pseudojet i
# and cross checks distance for other pseudojets back to i
# Note that nndist, near_neighbour, eta and phi are *Vectors*
function update_nearest_neighbour_crosscheck!(nndist, near_neighbour,
    i::Int, from::Int, to::Int, eta, phi, R2)
    new_nndist = R2
    new_nn = i
    @inbounds @simd for j in from:to
        delta2 = dist(i, j, eta, phi)
        if delta2 < new_nndist
            new_nn = j
            new_nndist = delta2
        end
        if delta2 < nndist[j]
            nndist[j] = delta2
            near_neighbour[j] = i
        end
    end
    nndist[i] = new_nndist
    near_neighbour[i] = nn;
end

```

129 Note there are some optimisations applied here as Julia *macros*, e.g., @simd, which we  
 130 discuss below. In particular, here we present updated results from those at the conference for  
 131 the tiled algorithm in Julia from applying the LoopVectorisation package in a key area  
 132 adding the @turbo macro:

```

find_best(diJ, n) = begin
    best = 1

```

```

134 @inbounds diJ_min = diJ[1]
135 @turbo for here in 2:n # Loop vectorisation marco
136     newmin = diJ[here] < diJ_min
137     best = newmin ? here : best
138     diJ_min = newmin ? diJ[here] : diJ_min
139 end
140 diJ_min, best
141 end

```

## 133 4 Code Performance

134 The different implementations of the anti- $k_T$  algorithm were tested on the same benchmark  
135 machine, a 64 core AMD EPYC 7302 3.00GHz with 24GB RAM, running CentOS7. The  
136 software versions used were gcc 11.3.0, Python 3.11.4 (with numba 0.57.1 and numpy 1.24.4)  
137 and Julia 1.9.2. More details on how to reproduce the measurements are given in [9].

138 Reconstruction of 100 LHC-like pp events<sup>2</sup> was run multiple times and the average re-  
139 construction time per event is given in Table 2. These numbers are normalised to the FastJet  
140 tiled algorithm performance (which is 324  $\mu$ s per event on the benchmark machine). Multiple  
141 repeats of the benchmark were done and jitter was observed to be extremely low, < 1%, so  
142 is not given. In these measurements the time to read the events (in HepMC3 format) and the  
143 JIT time for Julia and numba is excluded.

Implementation	Basic Algorithm	Tiled Algorithm
C++ (FastJet)	16.4	1.00
Python (Pure)	504	110
Python (Accelerated)	28.5	113
Julia	2.83	0.94

**Table 2.** Relative run times for the reconstruction of 100 13TeV pp events, normalised to the time for FastJet’s tiled algorithm. Results are stable and reproducible on the benchmark machine at < 1%.

144 We observe that the benchmark C++ FastJet code, with the tiled algorithm, is one of  
145 fastest implementations. The increase in performance for the tiled code, over the plain one,  
146 is significant with the events we used, confirming this is both an excellent algorithm and  
147 implementation for LHC p-p data.

148 As expected, the pure Python codes run very slowly in comparison. More surprisingly,  
149 the accelerated Python codes have quite poor performance as well. This is due to the fact  
150 that not all parts of the algorithm can be accelerated - bookkeeping operations still run in  
151 normal Python and become dominant in the overall runtime. This is particularly true of the  
152 tiled algorithm, which deliberately reduces the work to be done (which can be parallelised  
153 and accelerated) at the cost of more bookkeeping. This significantly hurts the accelerated  
154 implementation, which ends up slower than the basic accelerated implementation; it is not  
155 even faster than the pure Python tiled implementation code.

156 Our Julia code exceeds the performance of FastJet code. In the case of the tiled algorithm,  
157 as noted in Section 3.3, a *loop vectorisation* optimisation was applied to the search across all  
158  $d_{ij}$  to find the minimum value, which results in a 15% improved runtime on x86 architectures

---

<sup>2</sup>Hard QCD  $2 \rightarrow 2$  processes generated with Pythia8 at 13TeV, with a minimum transverse momentum of 20 GeV.

159 cf. without this macro<sup>3</sup>. In the case of the basic algorithm the Julia code uses a structure of ar-  
160 rays layout, which the compiler can highly optimise; additional benefit is gained from macros  
161 like `@simd`, which allow the compiler to apply further optimisations, gaining an additional  
162 5%.

163 There are some comments regarding these optimisations that should be made: the Julia  
164 compiler attempts to use SIMD instructions in any case; when using the `@simd` macro the  
165 developer is guaranteeing iterations are safe to reorder and to overlap, and that floating point  
166 operations can be reordered; `@turbo` also replaces some special functions with implementa-  
167 tions that can be vectorised better, but may be of lower accuracy. Use of these macros may  
168 lead to different numerical results so must be carefully validated (in our case we have checked  
169 that they are safe). One advantage in Julia is that these macros, as well as `@fastmath`, can  
170 be used and validated on a case-by-case basis (cf. the C++ compiler options such as `-O3`  
171 or `-ftree-vectorize`, which are applied per compilation unit, but more than likely are  
172 actually used globally). In addition, the JIT strategy of Julia and Python’s numba will auto-  
173 matically target the binary architecture of the machine being used, avoiding portability issues  
174 that can hamper C++ compiled binaries on different microarchitectures.

## 175 5 Conclusions

176 We have implemented the anti- $k_T$  algorithm in a number of different languages and examined  
177 code ergonomics as well as run time performance. The benchmark C++ code from FastJet  
178 is well written, but the hardest to reason on correctness, due to the nature of the language.  
179 Python is excellent for code logic and flexibility, but has a very poor run time performance;  
180 accelerating with numpy and numba unfortunately takes much of this advantage away, yet  
181 still fails to achieve a competitive run time. Julia performs extremely well, with an excellent  
182 ‘out of the box’ run time. The Julia compiler is able to find significant speed-ups and features  
183 like broadcast operators help to keep code clean and quick. Further, applying optimisations  
184 in Julia through the use of macros is extremely easy for the programmer to exploit and result  
185 in Julia having the best performance of all the codes that we tested.

186 It should be noted that the optimisations found by the Julia compiler could also be applied  
187 to the FastJet code to close the gap. However, the authors’ experience is that doing this in  
188 C++ is considerably more difficult.

189 Ergonomically, C++ is also the most difficult language to use, with no package manager,  
190 no built in profiler, and where templates and memory management remain tricky. The breadth  
191 of libraries in C++ is impressive, although managing dependencies is not easy. In Python the  
192 situation is far better, albeit that the package managers are not quite standardised (pip vs.  
193 conda/mamba). Profiling and debugging when accelerated code is used in Python (which  
194 is how Python is used in data intensive science) is not easy, but package support in Python  
195 is really excellent. In Julia the ecosystem is very well integrated, with a built in package  
196 manager and excellent reproducibility. Julia libraries are not as extensive as for C++ and  
197 Python, although the speed of development of new scientific libraries (which is Julia’s target  
198 community) is picking up quickly and most areas are covered (see the discussion in Eschel et  
199 al. [5]). Debugging and profiling in Julia are very well integrated.

200 We conclude that expanding the use of Julia in high energy physics would be very worth-  
201 while, given its excellent performance and ergonomics.

---

<sup>3</sup>On Apple’s M2Pro chip, the advantage for Julia is more significant, with the final Julia code running  $\times 1.45$  faster than FastJet for the tiled implementation cases



202 **References**

203 [1] J. Pisarski, *History and adoption of programming languages in NHEP* (2022),  
204 <https://indico.jlab.org/event/505/contributions/9207/>

205 [2] J. Bezanson, A. Edelman, S. Karpinski, V.B. Shah, *SIAM Review* **59**, 65 (2017)

206 [3] J. Bezanson, J. Chen, B. Chung, S. Karpinski, V.B. Shah, J. Vitek, L. Zoubritzky, *Proc.*  
207 *ACM Program. Lang.* **2** (2018)

208 [4] M. Stanitzki, J. Strube, *Comput. Softw. Big Sci.* **5**, 10 (2021), **2003**.11952

209 [5] J. Eschle, T. Gal, M. Giordano, P. Gras, B. Hegner, L. Heinrich, U.H. Acosta, S. Kluth,  
210 J. Ling, P. Mato et al., *Potential of the Julia programming language for high energy*  
211 *physics computing* (2023), **2306**.03675,  
212 <https://doi.org/10.48550/arXiv.2306.03675>

213 [6] J.M. Perkel, *Nature* **572**, 141 (2019)

214 [7] M. Miller, *Trends, challenges, and strategic shifts in the software vulnerability*  
215 *mitigation landscape* (2019), [https://github.com/Microsoft/](https://github.com/Microsoft/MSRC-Security-Research/blob/master/presentations/2019_02_BlueHatIL/2019_01%20-%20BlueHatIL%20-%20Trends%2C%20challenge%2C%20and%20shifts%20in%20software%20vulnerability%20mitigation.pdf)  
216 *MSRC-Security-Research/blob/master/presentations/2019\_02\_*  
217 *BlueHatIL/2019\_01%20-%20BlueHatIL%20-%20Trends%2C%20challenge%2C%*  
218 *20and%20shifts%20in%20software%20vulnerability%20mitigation.pdf*

219 [8] A. Naumann, *Journal of Physics: Conference Series* **513**, 052023 (2014)

220 [9] G.A. Stewart, P. Gras, B. Hegner, A. Krasnopolski, *Polyglot jet finding - LaTeX sources*  
221 *and benchmark instructions* (2023), <https://doi.org/10.5281/zenodo.8307668>

222 [10] M. Cacciari, G.P. Salam, *Phys. Lett. B* **641**, 57 (2006), hep-ph/0512210

223 [11] M. Cacciari, G.P. Salam, G. Soyez, *Journal of High Energy Physics* **2008**, 063 (2008)

224 [12] M. Cacciari, G.P. Salam, G. Soyez, *Eur. Phys. J. C* **72**, 1896 (2012), 1111.6097

225 [13] *Fastjet*, <https://fastjet.fr>