1 # CC-FJpy: A Python Package for seismic ambient noise

2 # cross-correlation and the frequency-Bessel transform method

3 **Authors:**

4 Zhengbo Li[1,2], Jie Zhou[5], Gaoxiong Wu[3], Jiannan Wang[3], Gongheng

5 Zhang[3], Sheng Dong[6], Lei Pan[3], Zhentao Yang[3], Lina Gao[3], Qingbo Ma[6],

6 Hengxin Ren[3], & Xiaofei Chen[2,3,4*]

7 1. Academy for Advanced Interdisciplinary Studies, Southern University of Science and
8 Technology, Shenzhen, 518055, China
9 2. Shenzhen Key Laboratory of Deep Offshore Oil and Gas Exploration Technology, Southern
10 University of Science and Technology, Shenzhen, 518055, China
11 3. Department of Earth and Space Sciences, Southern University of Science and Technology,
12 Shenzhen, 518055, China
13 4. Southern Marine Science and Engineering Guangdong Laboratory (Guangzhou), Guangzhou,
14 511458, China
15 5. School of Earth and Space Sciences, Peking University, Beijing, 100000, China
16 6. School of Earth and Space Sciences, University of Science and Technology of China, Heifei,
17 230026, China

18 ## Abstract

19 In the past two decades, surface wave imaging based on seismic ambient noise cross-correlation
20 (CC) has been one of the most important technologies in the field of seismology. With the
21 development of this technology, high-mode surface waves have received increasing attention,
22 especially after the proposition of the frequency-Bessel transform (F-J) method, which can
23 effectively extract multimode dispersion curves from ambient noise data. In the past few years,
24 our research group has made many attempts to improve this method. We summarized these
25 experiences and the corresponding algorithm for fast CC, and packaged them into a Python
26 package called CC-FJpy. It is commonly understood that CC takes a good deal of time. However,
27 we found that a simple reorganization of the CC logic can achieve computational acceleration by a
28 multiple of tens or even hundreds in comparison with classical CC open-source programs for $N$
29 stations. For the F-J method, we use Nvidia's graphics processing unit (GPU) to speed up
30 computation, and this approach achieves a hundreds-fold computational acceleration. We have
31 encapsulated our experiences and technologies into CC-FJpy and submitted it to various types of
32 data tests to ensure its speed and ease of use. We hope that providing the open source of CC-FJpy
33 can benefit the development of surface wave studies and make it easier to start with high-mode
34 surface waves. We look forward to your use and valuable suggestions.

35 ## Introduction

36 In the past two decades, significant understandings of underground structures of different

37  scales have been facilitated by the development of surface wave imaging with noise
38  cross-correlation (CC) technology (e.g., Campillo & Paul, 2003; Shapiro et al. 2005; Sabra et al.,
39  2005a, b; Yao et al. 2006; Bensen et al. 2009; Lin et al., 2009, 2011; Fang et al, 2015, 2016; Shen
40  et al., 2016). For ambient noise surface wave imaging, especially for lithospheric imaging, most
41  often the fundamental mode is0 obtained and inversed (e.g., Bensen 2007). Numerous studies
42  have confirmed that high-mode surface wave dispersion curves can provide more constraints on
43  underground structures (e.g., Nolet & Panza, 1976; Yokoi, 2010; Pan et al., 2018; Wu et al., 2020).
44  Wang et al. (2019a) proposed the frequency-Bessel transform (F-J) method, which can efficiently
45  extract Rayleigh wave multimode dispersion curves from ambient noise cross-correlation
46  functions (CCFs). Hu et al. (2020) verified that this method can be easily applied to Love waves;
47  Li & Chen (2020a; 2020b) extended this method to the application to seismic records, and
48  confirmed that this method can also extract the dispersion of PL waves. Zhan et al. (2020) applied
49  this method to imaging in Northeast China and updated the local 3-dimensional velocity model. To
50  further promote studies on high-mode surface waves and to ensure that more scholars can easily
51  use this method, we summarized our experiences in recent years and packaged our GPU F-J code
52  with our recently developed fast CC programs into an open-source Python package CC-FJpy.

53      Noise cross-correlation technology is one of the most important technologies in seismology.
54  CCFs obtained by CC can be approximated as Green's functions, which means that a large number
55  of seismological methods no longer rely on local earthquakes (e.g., Weaver & Lobkis 2004;
56  Sánchez-Sesma & Campillo 2006). CCFs have been widely used in surface wave imaging (e.g.,
57  Yao et al., 2006; Bensen et al., 2009), body wave imaging (e.g., Poli et al., 2012; Feng et al., 2017),
58  full wave inversion (Sager et al., 2017, 2020; Wang et al., 2019b), attenuation emulation (e.g.,
59  Lawrence et al., 2013) and so on. It is commonly understood that the CC process is often
60  time-consuming, especially when the overlap of time is needed (Seats et al., 2012). Ventosa et al.
61  (2019) attempted to accelerate the CC though GPUs. Although they accelerated the process of a
62  single CC for two stations, they could not accelerate the CC of $N$ stations well. After we carefully
63  studied the CC process and some widely used CC codes, we found that although CC technology
64  has been widely used for more than ten years since the early application of CC technology, there is
65  still a relatively large optimization space. For $N$ stations, $C_N^2$ times CCs are required. In many
66  classic programs, each CC between two stations comprises reading data, preprocessing and CC. In
67  fact, only $N$ reading data and preprocessing steps are required. Furthermore, the essence of CC
68  between records A and B is multiplication in the frequency domain:

$$CC = A(\omega) * conj(B(\omega)), \#1$$

69  where $conj$ is the conjugation. All the classic programs pack this step as a function (for example,
70  the MATLAB function $xcorr$), which means that every CC needs two fast Fourier transforms
71  (FFTs), which causes many repetitive FFT calculations. The total number of FFTs called is $2C_N^2$,
72  which can also be reduced to $N$. Based on these two points, we adjusted the logic of CC, wrote the
73  kernel in the C language and encapsulated it as a Python interface through Cython. Although this
74  sounds like a simple change, the effect is surprisingly good: its efficiency is ten to hundreds of
75  times higher than that of most classic CC programs. More importantly, our programs have very
76  small requirements for computing resources. In many cases, simple parallelism on a laptop is
77  enough to make it dozens of times faster than traditional programs on a server. In addition, our
78  program is also very easy to modify to adopt different kinds of improvements (e.g., Shen et al.,
79  2012; Xie et al., 2020) to the CC equation (equation 1).

80     For the F-J method, the core is to numerically realize the integral of equation 2.

$$I(\omega,k) = \int_0^\infty G(r,\omega)J_0(kr)r\,dr, \#2$$

81  where $k$ is the wavenumber, $r$ is the epicenter distance, $J_0(x)$ is the $0^{th}$ Bessel function of the first
82  kind and $G(r,\omega)$ can be CCFs or earthquake records. Generally, a trapezoidal integral is a good
83  choice, but this ignores the known characteristics of the Bessel function. Wang et al. (2019a) gave
84  a more accurate integration format, which will, however, increase the amount of calculation. To
85  balance efficiency and accuracy, we use the GPU to accelerate the process. The GPU is very
86  suitable for this type of calculation and can achieve hundreds of speedups, which can shorten the
87  original FJ process from a range of tens of minutes to hours to one of tens of seconds to minutes.
88  In addition, we encapsulated the integration of using the Hankel function instead of the Bessel
89  function, which has proven to be effective in removing "crossed" artifacts (Forbriger, 2003).
90  Different integration methods and GPU or non-GPU support are provided in CC-FJpy to facilitate
91  the needs of different users.

92     In recent decades, computer technology has brought revolutionary changes to many
93  industries. One of the most important drivers of these changes is that programming language and
94  complex algorithms have been efficiently encapsulated, so that numerous participants can quickly
95  learn and master the developed technology. The most typical example is the development of
96  machine learning. Now, even a middle school student can train his or her own model using
97  TensorFlow, PyTorch or other Python machine learning packages. In the field of geophysics, there
98  are also many informed scholars who have developed efficient open-source software programs,
99  such as the Generic Mapping Tools (GMT, https://www.generic-mapping-tools.org/) and Obspy
100 (https://docs.obspy.org/, Beyreuther et al. 2010). Encouraged by this, we decided to share our
101 small contribution in the direction of CC and high-mode surface waves to serve all colleagues, and
102 we have committed to maintaining the update for the foreseeable future. You can obtain CCFJpy
103 from https://github.com/ColinLii/CC-FJpy. We hope that through our program package, CC and
104 the extraction of high-order dispersion through the F-J method will become easier, especially for
105 scholars who are beginning to study this area. In addition, we humbly hope for valuable
106 suggestions.

107 **Implementation**

108     The imaging process through the F-J method can be simply summarized as the following four
109 steps: ① read data & preprocess, ② cross-correlation, ③ F-J scan and ④ dispersion curve
110 extraction & inversion (Figure 1). Among them, ① mainly depends on the storage format (e.g.,
111 SAC or miniSeed) and storage order of the data, and ④ has a lot of personalized solutions (e.g.,
112 Shen et al., 2012; Pan et al. 2018; Dereiling et al., 2019). We highly recommend Obspy for
113 reading and preprocessing data (https://docs.obspy.org/, Beyreuther et al. 2010). The calculations
114 of ② and ③ are relatively fixed. Thus, CC-FJpy mainly deals with ② and ③ and is divided
115 into two sub-packages: CCpy and FJpy. The two sub-packages can be used together or completely
116 independently. In addition, although ④ has a large number of personalized programs, we plan to
117 add several inversion methods that we believe are efficient and robust as examples in future
118 updates.

119 **CCPY: a Python sub-package for rapid cross-correlation**
120     First, let us briefly review the basic formula of cross-correlation:

$$C_{1,2}(t) \approx \int_0^{t_C} v_1(\tau)v_2(t+\tau)d\tau, \#3$$

121 where $v_1(t)$ and $v_2(t)$ are the continuous broadband records of stations 1 and 2 in the time
122 window $[0, t_C]$. Usually, $t_C$ is not the total continuous recording time, as the recording is divided
123 into unit time lengths such as one hour, one day or one week. CCFs are obtained by superimposing
124 a large number of $C_{1,2}(t)$ with different time windows. The realization of equation 3 in the time
125 domain is time-consuming, so for most programs, it is implemented in the frequency domain.

$$C_{1,2}(\omega) \approx dft(v_1)conj(dft(v_2)), \#4$$

126 where $dft$ is the discontinuous Fourier transform and $conj$ is the conjugation.

127 From equation 3 and equation 4, the CC of the two stations is appears to be very simple and
128 without much room for acceleration. However, the strategy for calculating the CC affects the
129 calculation time. We first talk about the outputs as the frequency-domain CC functions $C_{j,k}(\omega)$
130 for $N$ stations. The classic strategy is what we called strategy 1, the completely independent
131 strategy, which means that every time the two stations are cross-correlated, the data of the two
132 stations are read and correlated (Figure 2a). For this strategy, for one CC time unit, $N(N\text{-}1)$ times
133 reading data and $N(N\text{-}1)/2$ times CC are required. Obviously, there are many duplications in the
134 reading data step. An improvement in this strategy is to read the data of $N$ stations in once and
135 then $N(N\text{-}1)/2$ times CC can be performed (Figure 2b). We call this strategy 2, the shared memory
136 strategy. Compared with the completely independent strategy, the number of readings drops from
137 $N(N\text{-}1)$ times to N times. Furthermore, according to equation 4, we can divide CC into two parts:
138 the FFT and multiplication (Figure 2d). The FFT can be shared like the reading data. Thus, we
139 have strategy 3: shared memory and the FFT strategy (Figure 2c). In addition to the reading time,
140 the CC time can also be reduced. It is worth mentioning that we use the C language to call the
141 fftw-3 package (https://www.fftw.org) for FFT, as it is approximately 3 times faster than Python
142 numpy fft. All the C codes are encapsulated as a Python interface through Cython.

143 For the outputs are the time-domain CC functions $C_{i,j}(t)$, the shared memory and FFT
144 strategy have more significant time advantages for $M$ CC time units. For the first two strategies,
145 since each CC outputs the time-domain CC functions, $M \times N \times (N-1)/2$ inverse FFTs (IFFTs)
146 are needed before the overlap. However, since the Fourier transform is linear, we can overlap in
147 the frequency domain and then perform the IFFT, which means that only $N \times (N-1)/2$ IFFTs
148 are needed. According to our test, the multiplication takes much less time than reading and
149 preprocessing data, the FFT and the IFFT. The efficiency of the acceleration through strategy 3 is
150 positively related to the total number of stations $N$ and the total number of stacking days $M$. The
151 larger $N$ and $M$ are, the more obvious the improvement delivered by strategy 3. Different machines
152 and different data will have a great impact on the time cost. After using strategy 3, the overall CC
153 efficiency is improved by one to two orders over strategies 1 and 2. As the amount of calculation
154 is greatly reduced, CCpy is suitable for both personal PCs and notebooks.

155 In many cases, to improve the quality of the CC, overlaps of units are needed (Seats et al.,
156 2012). In our program, we designed a larger reading unit $T_C$ for reading, and the overlap of $t_C$ is
157 executed within to improve efficiency. You can select an overlap rating, whether to use spectral
158 whitening, whether to use onebit and other options for different situations with the interface
159 ccfj.cc. We will show you a specific example of data from USArray in the next section to show the
160 acceleration efficiency. For details, please read the package manual.
161 **FJPY: a Python sub-package for the F-J method through GPU**

162    The core of the F-J method is the numerical realization of equation 2. For $N$ observed $G(r_i,\omega)$

163    arranged in order from station distance for CCFs or epicenter distance for earthquake records, the

164    trapezoidal integral can be used to approximate equation 2:

$$I(\omega,k) \approx \frac{1}{2}\sum_{i=1}^{N-1}(G(r_i,\omega)J_0(kr_i) + G(r_{i+1},\omega)J_0(kr_{i+1}))(r_{i+1} - r_i).\#5$$

165    Note that for trapezoidal integration, $G(r,\omega)$ is only obtained at the observation points $r_i$, but

166    $J_0(kr)$ is known from 0 to $\infty$. Thus, Wang et al. (2019a) gave another numerical integral format

167    of equation 4 through linear approximation of Green's function:

$$I(\omega,k) \approx \sum_{i=1}^{N-1}\left\{\frac{1}{k}G(r,\omega)rJ_1(kr) + \frac{b_j}{k^3}[krJ_0(kr) - B_0(kr)]\right\}\Big|_{r_j}^{r_{j+1}},\#6$$

168    where $b_j = \frac{G(r_{j+1},\omega)-G(r_j,\omega)}{r_{j+1}-r_j}$, and $B_0(x) = \int_0^x J_0(\eta)d\eta$. In the early implementation of the F-J

169    method, the calculation of $B_0(x)$ is by trapezoidal integration. Later, we find the primitive of

170    $B_0(x)$:

$$\int J_0(x)dx = xJ_0(x) + \frac{\pi x}{2}[J_1(x)Q_0(x) - J_0(x)Q_1(x)],\#7$$

171    where $Q_i(x)$ is the i[th] Struve function, which can be calculated by the subroutine of Ruckdeschel

172    (1981).

173        For dispersion features, seismologists prefer to display in the frequency-phase-velocity (f-c)

174    domain over the frequency-wavenumber (f-k) domain. The domain conversion can be performed

175    by:

$$c = \frac{2\pi\omega}{k}.\#8$$

176    Commonly, when calculating *nc* phase velocity points and *nf* frequency points, regardless of

177    whether equation 5 or equation 6 is used, the size of the calculation of the F-J method is quite

178    large, especially for noise data. However, the F-J method is naturally suitable for parallel

179    acceleration through the Nvidia GPU. Each calculation of $c_i$ and $\omega_i$ is not related to each other.

180    Compute unified device architecture (CUDA) programming enables us to execute F-J integration

181    on GPU devices. Equation 5 and 6 can be packaged into different 'kernels' which is the code run

182    on the GPU device, and be scheduled by *nf*×*nc* GPU threads. Further, we encapsulate the CUDA

183    program with Python; this makes it possible to quickly implement equations 5 and 6 by calling

184    function ccfj.fj with different parameters. It is worth noting that changing the Bessel function in

185    equations 5 and 6 into the first kind of Hankel function will help eliminate the "cross" artifact

186    (Forbriger, 2003). We have also added parameters to control using the Bessel function or Hankel

187    function.

188        It is also worth noting that for noise data, we often only use the real part of the CCFs for

189    calculation, while for seismic data, we calculate the real and imaginary parts of the recorded

190    spectrum and take $|I(\omega,k)|$. This is mainly because an earthquake has a source time function,

191    which will affect the results of the pure real or imaginary part. In addition, Li & Chen (2020)

192    noted that the F-J method for seismic records often requires auxiliary time windows

193    (multi-windows F-J method, MWFJ). Therefore, we specifically designed the ccfj.mwfj interface

194    for earthquake events. For details, please refer to the manual and the Python examples.

## Examples of USArray

Two application examples are illustrated: one is an ambient noise example consistent with Wu et al. (2020) and the other is an earthquake example consistent with Li & Chen (2020).

**Ambient Noise**

The data we use are half a year (182) of continuous records from June $1^{st}$ (day 152) to December $1^{st}$ (day 334) in 2011 of 96 stations from the USArray (Figure 3a). The original data size is approximately 55GB. We cropped data by day ($T_C$), which is probably the most common split time. Then, we downsampled to 4 Hz, demeaned, detrended, removed the instrument response and saved it as SAC files. After decompression, the size of a single SAC file is 1.31 MB and the total file size is 20.9 GB. If the data are stored as a longer period ($T_C$), such as week or month, the reading efficiency and calculation efficiency will be higher. Hourly CC with 90% time overlaps and spectral whitening is adopted during CC. Figure 3b shows the CCFs in the frequency-domain recovered by CC-FJpy, while Figure 3c shows the time-domain CCFs obtained by the IFFT of the frequency CCFs. Both the frequency-domain and time-domain CCFs have good coherence.

As mentioned in the last section, we care most about the computational efficiency. Please note that different machines and data will have a greater impact on the results. The CPU applied for the test was a 10-cores Intel(R) Core (TM) i9-10900K with 64 GB ddr4 2666 MHz RAM and Seagate Exos 7E8 ROM. For the accuracy of the test, we read 100 different SAC files of one day and performed demeaning and detrending. The average time of reading data, demeaning and detrending is approximately 0.05 seconds. Similarly, we calculated that the time required to calculate the FFT with numpy is 0.093 seconds while that with fftw-3 is 0.029 seconds for hourly CC with 90% overlap at one station. The multiplication time of the numpy array is 0.008 seconds while the of fftw-3 is 0.002 seconds (Figure 4a). Figure 4b shows that the different strategies need to calculate the number of times to read data & demean & detrend, conduct the FFT and multiplication. The number of reads and FFTs required by strategy3 has drops sharply compared to the other two strategies, which leads to the CC time for 96 stations in a day being much less than that of other two strategies. To calculate CC in the frequency-domain for 96 stations in one day, strategy 1 takes 815 seconds, strategy 2 takes 542 seconds, and strategy 3 takes 11.5 seconds. Figure 4c-e shows the percentage of the different strategies, where the area is proportional to the time used, and the read and FFT time saved by strategy 3 is easily seen. For CCpy, which uses strategy 3, it takes less than 1800 seconds to complete the CC of the 4 Hz data of 96 stations for half a year in series. Under parallelism, since a large part of the time in the cross-correlation is reading data, the efficiency of parallelism does not entirely depend on the number of cores, but also depends on the speed of the hard disk. It takes less than 10 minutes to use 20 threads in parallel. For the first two strategies, it takes more than 24 hours to use serial and at least 4 hours to use parallel. It is worth noting that the time is measured on the author's personal computer, which may be quite unstable.

It should be noted that the above comparison does not consider the time taken by the IFFT. If the time of the IFFT is considered, as Figure 2d shows, for strategies 1 and 2, every CC needs a 1-time IFFT, which means that for one-day data with 90% overlap, $231 \times C_N^2$ IFFTs are needed, and for 182-day data, $182 \times 231 \times C_N^2$ IFFTs are needed. However, for strategy 3, only $C_N^2$ IFFTs are required, which will further highlight the acceleration ratio of strategy 3 compared to

238 those of strategies 1and 2.

239      We use the trapezoidal integral (equation 5) based on the Bessel function, the linear
240 approximate integral (equation 6) based on the Bessel function, the trapezoidal integral (equation
241 5) based on the Hankel function and the linear approximate integral (equation 6) based on the
242 Hankel function to extract the dispersion spectrum from the CCFs (Figure 5). Compared with
243 trapezoidal integration, linear approximate integration can improve the quality of the dispersion
244 spectrum. The Hankel function can effectively remove "cross" artifacts. Since GPU acceleration is
245 used, the calculation time is approximately tens of seconds, and the specific values of the time cost
246 are marked in the subfigures in Figure 5. The GPU applied in the test is the Nvidia RTX 2070
247 super, which is commonly used and at a suitable price. Both the Linux platform and the Windows
248 platform are supported. We also provide the corresponding calculation program without GPU
249 acceleration in the program package, but the calculation without a GPU is relatively slow. If it is
250 completely serialized, it will take close to 3 hours of calculation for the linear approximate integral
251 of the Bessel function (Figure 5b), while the calculation with GPU acceleration is approximately
252 11 seconds. Therefore, we strongly recommend using an Nvidia GPU to accelerate; even a very
253 ordinary GPU will achieve a high acceleration.

254 **Earthquake**

255      Here, we repeat the example of the Mw 5.7 Oklahoma earthquake in Li & Chen (2020) with
256 the FJPY. In this case, MWFJ with three time windows, NoWin, which means no time window,
257 Win1 [3.2, 3.7] km/s and Win2 [3.7, 4.3] km/s, are applied. Here, we only show the results
258 calculated according to the Bessel function and Hankel function corresponding to equation 6.
259 Figures a, b and c show the dispersion spectrum extracted with three time windows by equation 6
260 with the Bessel function, while Figures d, e and f show the dispersion spectrum extracted with
261 three time windows by equation 6 with the Hankel function. The specific code calls have been
262 shown in "earthquake.ipynb". We prefer users try both instead of comparing the results, as, the
263 calculation of seismic records is generally approximately a few seconds.

# Discussion and Conclusions

265      The F-J method is an effective method for extracting high-mode surface wave dispersion
266 from various types of seismic records, and has recently received increasing attention. We
267 summarized the application of our group's research in recent years on the F-J method, and
268 encapsulated our codes and experiences into a Python package. We hope that through open source,
269 we have made it more convenient for more seismologists to use the F-J method.

270      At present, this package contains two parts: CCpy, which performs fast noise
271 cross-correlation, and FJpy, which is accelerated by Nvidia's GPU. Although we only made minor
272 modifications to the existing cross-correlation logic, CCpy still delivers several times more speed
273 up, so that CC that used to take days or weeks will now only take tens of minutes to a few hours.
274 We believe this is helpful not only for the F-J method with ambient noise but also for many other
275 seismological studies. The GPU acceleration drops the time taken by the F-J method, especially
276 the application of ambient noise, from tens of minutes to approximately 1 minute, which greatly
277 improves the efficiency of F-J imaging. The GPU can also be used to accelerate the CC process
278 (Ventosa et al., 2019), and we are considering adding it in a future update. However, we are
279 concerned about that GPUs cannot bring qualitative acceleration, such as CCpy, because reading
280 data is required. As shown in the example in Figure 4, even if the FFT time and the multiplication

281 time are both 0, the reading still needs more than 1/3 of the time. However, with the increasing
282 popularity of dense arrays, especially the application of distributed acoustic sensing (DAS,
283 Mateeva et al., 2014; Hartog, 2017) technology, any attempts to improve efficiency should be
284 encouraged.

285      We believe that with the popularity of computers today, the ease of use of codes is very
286 important for industry development. We chose to encapsulate our code in the form of a Python
287 package, which is very suitable for embedding existing codes and applications. However, we
288 believe that we are still inexperienced in developing and maintaining open-source code, so we will
289 continue to humbly seek valuable advice.

## Data and Resources

291 All the seismic records used in the examples were requested from the Data Management Center
292 (DMC) of Incorporated Research Institutions for Seismology (IRIS) at
293 https://ds.iris.edu/ds/nodes/dmc. Additionally, USArray information can be obtained from
294 https://www.usarray.org and https://doi.org/10.7914/SN/TA. Detailed information about fftw-3 can
295 be obtained at https://www.fftw.org. The CC-FJpy, manual and examples are available from
296 https://github.com/ColinLii/CC-FJpy. We have also uploaded the Jupyter notebook files of
297 examples in the supplemental materials. All the links mentioned are last accessed on February 4,
298 2021.

## Acknowledgements

309

## References

311 Bensen, G. D., Ritzwoller, M. H., Barmin, M. P., Levshin, A. L., Lin, F., Moschetti, M. P., . . .
312 Yang, Y. (2007). Processing seismic ambient noise data to obtain reliable broad-band surface wave
313 dispersion measurements. Geophysical Journal International, 169(3), 1239-1260.
314 doi:10.1111/j.1365-246X.2007.03374.x

315

316 Bensen, G. D., Ritzwoller, M. H., & Yang, Y. (2009). A 3-D shear velocity model of the crust and
317 uppermost mantle beneath the United States from ambient seismic noise. Geophysical Journal
318 International, 177(3), 1177-1196. doi:10.1111/j.1365-246X.2009.04125.x

319

320 Beyreuther, M., Barsch, R., Krischer, L., Megies, T., Behr, Y., & Wassermann, J. (2010). ObsPy: A
321 Python Toolbox for Seismology. Seismological Research Letters, 81(3), 530-533.

322    doi:10.1785/gssrl.81.3.530

324    Campillo, M., & Paul, A. (2003). Long-range correlations in the diffuse seismic coda. Science,
325    299(5606), 547-549. doi:10.1126/science.1078551

327    Dreiling, Jennifer; Tilmann, Frederik (2019): BayHunter - McMC transdimensional Bayesian
328    inversion of receiver functions and surface wave dispersion. GFZ Data Services.
329    http://doi.org/10.5880/GFZ.2.4.2019.001

331    Fang, H., Yao, H., Zhang, H., Huang, Y.-C., & van der Hilst, R. D. (2015). Direct inversion of
332    surface wave dispersion for three-dimensional shallow crustal structure based on ray tracing:
333    methodology and application. Geophysical Journal International, 201(3), 1251-1263.

335    Fang, H., Zhang, H., Yao, H., Allam, A., Zigone, D., Ben-Zion, Y., . . . van der Hilst, R. D. (2016).
336    A new algorithm for three-dimensional joint inversion of body wave and surface wave data and its
337    application to the Southern California plate boundary region. Journal of Geophysical Research:
338    Solid Earth, 121(5), 3557-3569.

340    Feng, J., Yao, H., Poli, P., Fang, L., Wu, Y., & Zhang, P. (2017). Depth variations of 410 km and
341    660 km discontinuities in eastern North China Craton revealed by ambient noise interferometry.
342    Geophysical Research Letters, 44(16), 8328-8335. doi:https://doi.org/10.1002/2017GL074263

344    Forbriger, T. (2003). Inversion of shallow-seismic wavefields: I. Wavefield transformation.
345    Geophysical Journal International, 153(3), 719-734. doi:10.1046/j.1365-246X.2003.01929.x

347    Hartog, A. H. (2017). An introduction to distributed optical fibre sensors. Boca Raton, Florida:
348    CRC Press/Taylor and Francis. https://doi.org/10.1201/9781315119014

350    Hu, S., Luo, S., & Yao, H. (2020). The Frequency-Bessel Spectrograms of Multicomponent
351    Cross-Correlation Functions From Seismic Ambient Noise. Journal of Geophysical Research:
352    Solid Earth, 125(8), e2020JB019630. doi:https://doi.org/10.1029/2020JB019630

354    Lawrence, J. F., Denolle, M., Seats, K. J., & Prieto, G. A. (2013). A numeric evaluation of
355    attenuation from ambient noise correlation functions. Journal of Geophysical Research: Solid
356    Earth, 118(12), 6134-6145. doi:https://doi.org/10.1002/2012JB009513

358    Li, Z., & Chen, X. (2020). An Effective Method to Extract Overtones of Surface Wave From Array
359    Seismic Records of Earthquake Events. Journal of Geophysical Research: Solid Earth, 125(3),
360    e2019JB018511. doi:10.1029/2019jb018511

362    Li, Z., & Chen, X. (2020). Multiple dispersion curves extracted from seismic PL phase. Earth and
363    Space Science Open Archive. doi: 10.1002/essoar.10505116.1

365    Mateeva, A., Lopez, J., Potters, H., Mestayer, J., Cox, B., Kiyashchenko, D., . . . Detomo, R.

366    (2014). Distributed acoustic sensing for reservoir monitoring with vertical seismic profiling.
367    Geophysical Prospecting, 62(4), 679-692. doi:https://doi.org/10.1111/1365-2478.12116

369    Nolet, A. M. H. (1976). Higher modes and the determination of upper mantle structure: Drukkerij
370    Elinkwijk.

372    Pan, L., Chen, X., Wang, J., Yang, Z., & Zhang, D. (2018). Sensitivity analysis of dispersion
373    curves of Rayleigh waves with fundamental and higher modes. Geophysical Journal International,
374    216(2), 1276-1303.

376    Poli, P., Campillo, M., Pedersen, H., & Group, L. W. (2012). Body-wave imaging of Earth's
377    mantle discontinuities from ambient seismic noise. Science, 338(6110), 1063-1065. Retrieved
378    from https://science.sciencemag.org/content/sci/338/6110/1063.full.pdf

380    Ruckdeschel F. R. (1981) BASIC Scientific Subroutines, Vol. II.

382    Sabra, K. G., Gerstoft, P., Roux, P., Kuperman, W. A., & Fehler, M. C. (2005). Extracting
383    time-domain Green's function estimates from ambient seismic noise. Geophysical Research
384    Letters, 32(3).

386    Sabra, K. G., Gerstoft, P., Roux, P., Kuperman, W. A., & Fehler, M. C. (2005). Surface wave
387    tomography from microseisms in Southern California. Geophysical Research Letters, 32(14).

389    Sager, K., Boehm, C., Ermert, L., Krischer, L., & Fichtner, A. (2020). Global-Scale
390    Full-Waveform Ambient Noise Inversion. Journal of Geophysical Research: Solid Earth, 125(4),
391    e2019JB018644. doi:https://doi.org/10.1029/2019JB018644

393    Sager, K., Ermert, L., Boehm, C., & Fichtner, A. (2017). Towards full waveform ambient noise
394    inversion. Geophysical Journal International, 212(1), 566-590. doi:10.1093/gji/ggx429

396    Sánchez-Sesma, F. J., & Campillo, M. (2006). Retrieval of the Green's function from cross
397    correlation: the canonical elastic problem. Bulletin of the Seismological Society of America, 96(3),
398    1182-1191. doi:10.1785/0120050181

400    Shapiro, N. M., Campillo, M., Stehly, L., & Ritzwoller, M. H. (2005). High-resolution
401    surface-wave tomography from ambient seismic noise. Science, 307(5715), 1615-1618.
402    doi:10.1126/science.1108339

404    Seats, K. J., Lawrence, J. F., & Prieto, G. A. (2012). Improved ambient noise correlation functions
405    using Welch′s method. Geophysical Journal International, 188(2), 513-523.
406    doi:10.1111/j.1365-246X.2011.05263.x

408    Shen, W., & Ritzwoller, M. H. (2016). Crustal and uppermost mantle structure beneath the United
409    States. Journal of Geophysical Research: Solid Earth, 121(6), 4306-4342.

410    doi:10.1002/2016JB012887

411

412    Shen, W., Ritzwoller, M. H., Schulte-Pelkum, V., & Lin, F.-C. (2012). Joint inversion of surface
413    wave dispersion and receiver functions: a Bayesian Monte-Carlo approach. Geophysical Journal
414    International, 192(2), 807-836.

415

416    Sergi Ventosa, Martin Schimmel, Eleonore Stutzmann; Towards the Processing of Large Data
417    Volumes with Phase Cross-Correlation. Seismological Research Letters 2019;; 90 (4): 1663–1669.
418    doi: https://doi.org/10.1785/0220190022

419

420    Wang, J., Wu, G., & Chen, X. (2019). Frequency-Bessel Transform Method for Effective Imaging
421    of Higher-Mode Rayleigh Dispersion Curves From Ambient Seismic Noise Data. Journal of
422    Geophysical Research: Solid Earth, 124(4), 3708-3723. doi:10.1029/2018jb016595

423

424    Wang, K., Liu, Q., & Yang, Y. (2019). Three-Dimensional Sensitivity Kernels for Multicomponent
425    Empirical Green's Functions From Ambient Noise: Methodology and Application to Adjoint
426    Tomography. Journal of Geophysical Research: Solid Earth, 124(6), 5794-5810.
427    doi:https://doi.org/10.1029/2018JB017020

428

429    Weaver, R. L., & Lobkis, O. I. (2004). Diffuse fields in open systems and the emergence of the
430    Green's function (L). The Journal of the Acoustical Society of America, 116(5), 2731-2734.
431    doi:10.1121/1.1810232

432

433    Wu, G.-x., Pan, L., Wang, J.-n., & Chen, X. (2020). Shear Velocity Inversion Using Multimodal
434    Dispersion Curves From Ambient Seismic Noise Data of USArray Transportable Array. Journal of
435    Geophysical Research: Solid Earth, 125(1), e2019JB018213. doi:10.1029/2019jb018213

436

437    Yao, H., van Der Hilst, R. D., & De Hoop, M. V. (2006). Surface-wave array tomography in SE
438    Tibet from ambient seismic noise and two-station analysis—I. Phase velocity maps. Geophysical
439    Journal International, 166(2), 732-744. doi:10.1111/j.1365-246X.2006.03028.x

440

441    Yokoi, T. (2010). New formulas derived from seismic interferometry to simulate phase velocity
442    estimates from correlation methods using microtremorFormulas for phase velocity estimates.
443    Geophysics, 75(4), SA71-SA83.

444

445    Zhan W., Pan L., Chen X.F., (2020). "3D crustal shear velocity structure of China Northeast
446    constrained from inversion of multimodal dispersion curves of Rayleigh waves from ambient
447    seismic noise", J. Asian Earth Sci., 196, July, doi.org/10.1016/j.jseaes.2020.104372

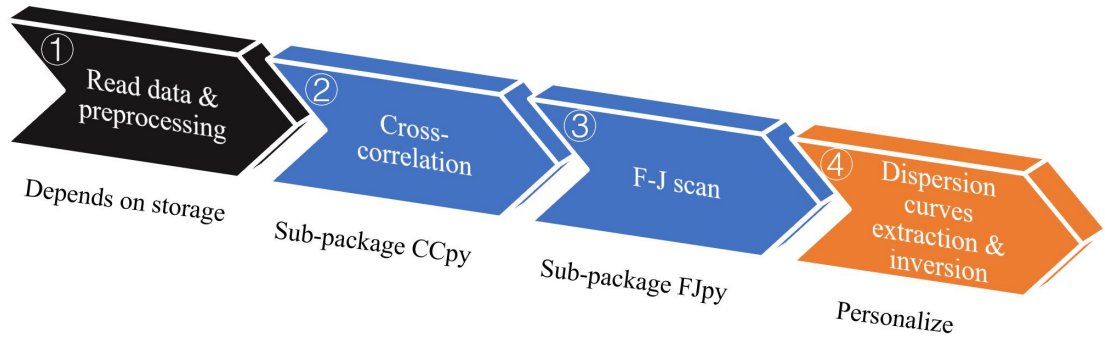448

449    **Figures**

450

451    **Figure 1.** The imaging process through the F-J method. At present, CC-FJpy mainly contains

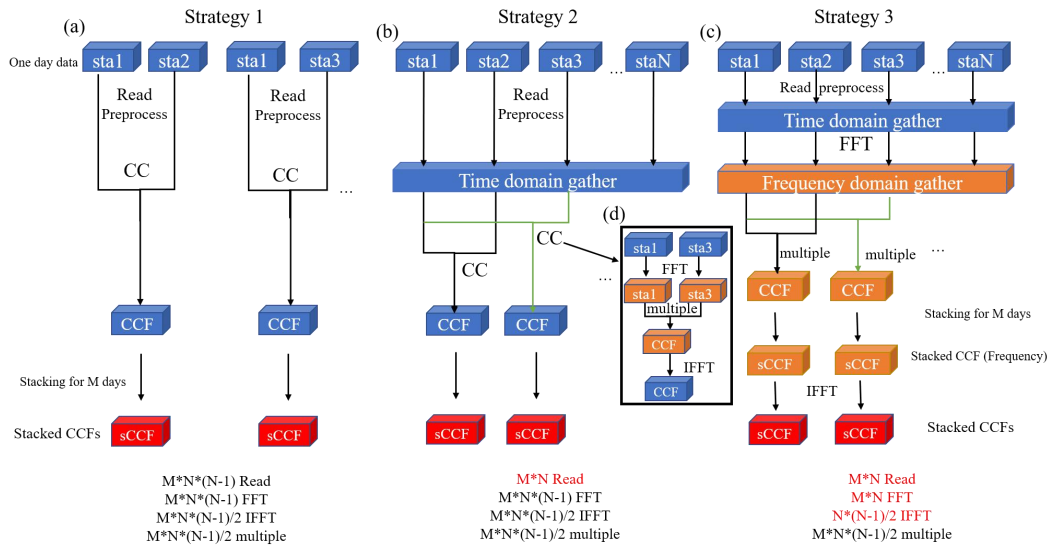452    parts ② and ③. We will add ④ to this package in future updates



453

454    **Figure 2.** Comparison of the three CC strategies for the continuous records of $M$ days and $N$

455    stations (sta1, sta2, …, staN) in units of days. The blue squares are the data in the time domain, the

456    yellow squares are the data in the frequency domain, and the red squares are the final outputs. **(a)**

457    Completely independent cross-correlation: in this approach, each CC reads data and correlates

458    independently. **(b)** Shared memory strategy: in this approach, $N$ records of each day are read,

459    shared and then cross-correlated. **(c)** Shared reading and FFT strategy: not only are the memories

460    of N records shared the FFTs and IFFTs are also shared. **(d)** How a single CC is implemented.

461    There is a huge gap in the required calculation between the three strategies. We have marked the

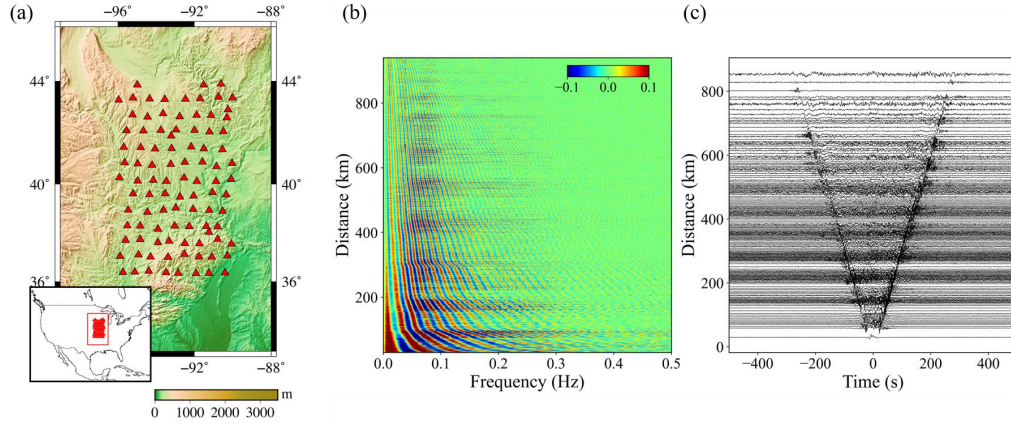462    number of core calculations for CC below the flowchart of each strategy.

463

464

465 **Figure 3.** The CCFs recovered from USArray ambient noise data. (a) Stations used. (b) CCFs in
466 the frequency domain recovered by CC-FJpy. (c) Time domain CCFs obtained by the IFFT of the
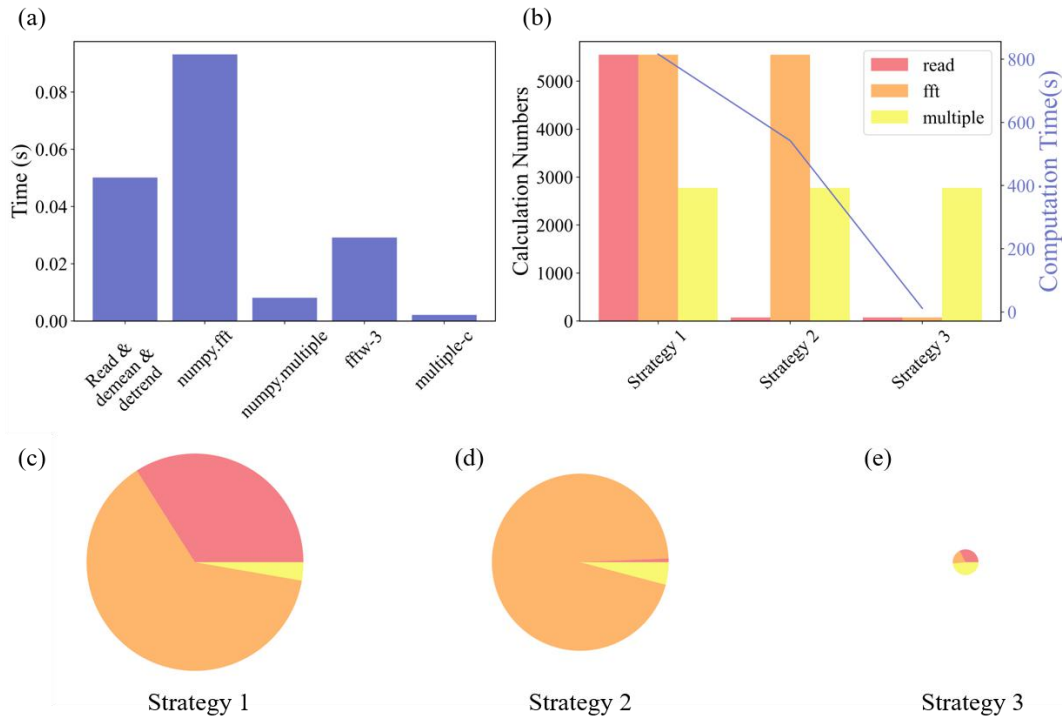467 frequency-domain CCFs.

468



469

470 **Figure 4.** Cross-correlation performance comparison. It should be noted that all the
471 comparisons here are in the case of a serial approach. **(a)** The time taken to read and preprocess
472 data, FFT by numpy, multiplication by Python, FFT by fftw-3 and multiplication by C language
473 for one station per day data. It should be noted that the FFT times and multiple times refer to the
474 sum time of FFT and multiple in the overlap. **(b)** Comparison of the calculation amounts of
475 different CC strategies. **(c), (d), and (e)** Proportion of time consumed by the main operations for
476 the different strategies. The area of the pie chart is proportional to the total time.
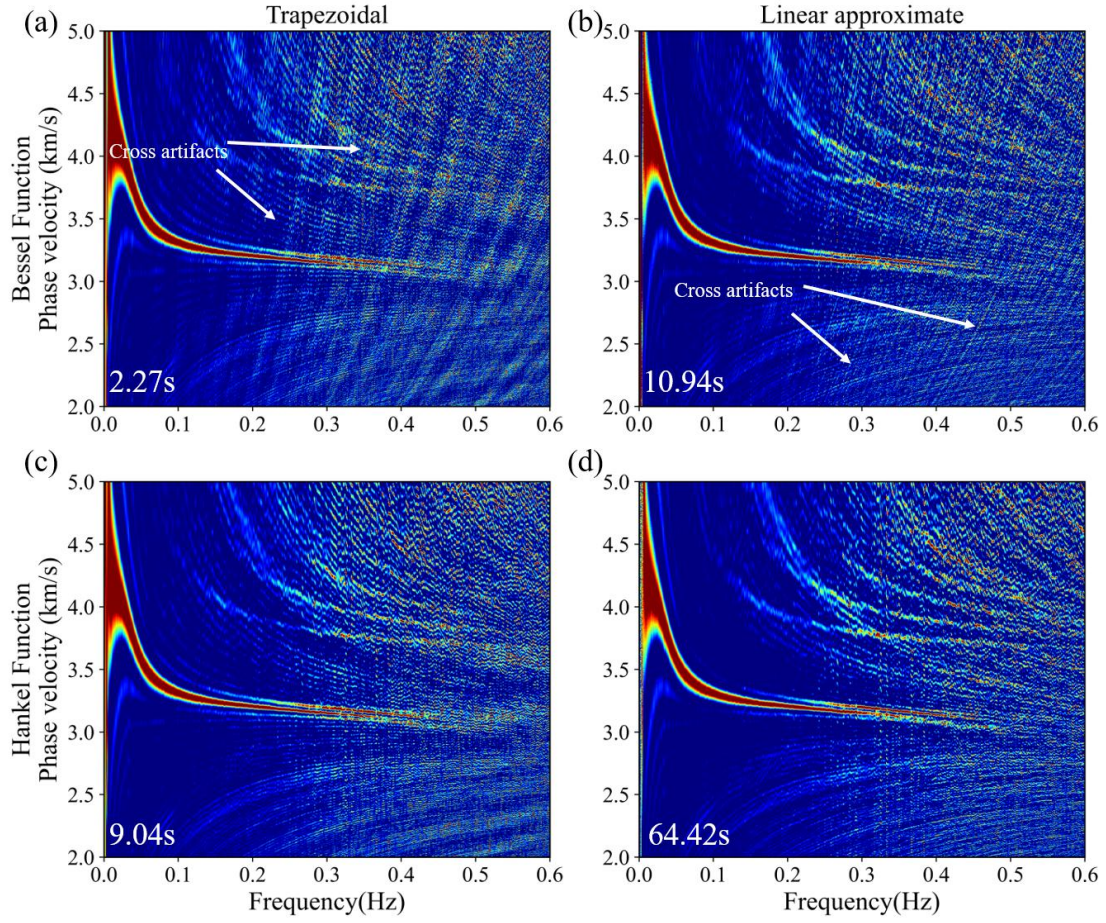
477

**Figure 5.** Dispersion spectra extracted from CCFs by FJpy. **(a)** Dispersion spectrum calculated through the trapezoidal integral of the Bessel function. **(b)** Dispersion spectrum calculated through equation 6. **(c)** Dispersion spectrum calculated through the trapezoidal integral of the Hankel function. **(d)** Dispersion spectrum calculated though equation 6 with the Hankel function. The times in the lower left corner of each subfigure are the calculation times after acceleration by the GPU.
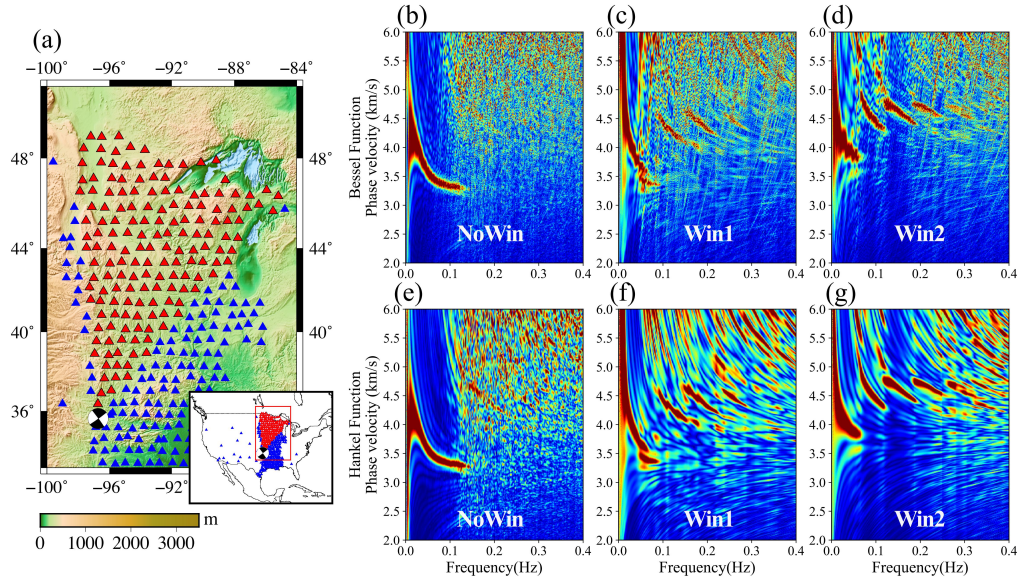
486
487 **Figure 6.** Dispersion spectra extracted from seismic records by FJpy. **(a), (b), and (c)** Dispersion
488 spectra extracted by equation 6 with NoWin, Win1 and Win2. **(d), (e), and (f)** Dispersion spectra
489 extracted by equation 6 with the Hankel function.
490