

# PFTSEARCH.DOC

(last update: November 14, 2003)

## CONTENTS

- A. INTRODUCTION
- B. GENERAL DESCRIPTION OF PROGRAM
- C. PFTSEARCH PROCEDURE
- D. PROGRAM OUTPUT
- E. BATCH JOB PROGRAM INPUT
- F. STRATEGY FOR USE OF PFTSEARCH AND EM3DR
- G. DISCLAIMER
- H. CREDITS
- I. FINAL NOTES
- J. PROGRAM MODIFICATIONS
- K. FLOW CHART FOR PFTSEARCH PROGRAM

## A. INTRODUCTION

PFTSEARCH determines orientations (THETA, PHI, OMEGA =  $\theta, \phi, \omega$ ) and origin (X, Y) values for particle images. This, of course, is critically important information for producing a 3D reconstructed density map from many, independent images. PFTSEARCH uses cross-correlation procedures to compare unfiltered or filtered particle images against a database of reference projection views produced from a 3D density map. The map might be a 3D reconstruction or model, for example, computed from an atomic, X-ray crystallographic structure. The last major enhancements to the program included use of compressed data and the sliding window approach to refinement (released in May 1998) and most recently the use of contrast transfer function (CTF) corrections as (August 1999).

### **Use of CTF to assist comparison of image and projection data**

The user has three different options for imposing the microscope CTF on the model projections or raw images with the goal of getting better refinement of orientation and origin parameters. The user, of course, has the option (set MODE2 = 0) of NOT using the CTF correction software.

The four CTF options are:

MODE2 = 1 Multiply the projections of the 3D model by the CTF.

MODE2 = 2 Multiply the raw image data by the CTF.

MODE2 = 3 Multiply the raw image data by the inverse CTF.

MODE2 = 1 is designed for the situation in which you have imposed a CTF correction in the 3D map in the program EM3DR. In theory, PFTSEARCH will work best if the model projections and raw

images are compared 'apples to apples'. That is to say, by imposing a CTF onto the model projections, you are trying to recreate what noise-free data (the model) would look like if the projection pictures had been recorded in the TEM like the raw particle images. **IMPORTANT WARNING:** In PFTSEARCH's present state use of this option is ONLY CORRECT for image data from ONE MICROGRAPH. If you are refining image data from multiple micrographs (with different CTFs), then the use of MODE2 = 1 \*\*\*may\*\*\* lead to errors in refinement.

MODE2 = 2 employs the strategy used by Steve Fuller's group. Here, raw images are multiplied by the microscope CTF and this has the effect of weighting the image data to emphasize the most significant (i.e. highest S/N) and de-emphasize the least significant (low S/N) portions of each image. This down weights data near the nodes of the CTF, precisely where the noise is highest in the image transforms. So, option 2 weights the data to emphasize the reliable portions of the image while down weighting the most unreliable portions. This, in turn, is designed to help the PFT routine perform a more reliable refinement of orientation and origin parameters.

MODE2 = 3 causes each raw image to be multiplied by the INVERSE CTF as a means to create image data that more properly compares with the UNMODIFIED model projections.

Only careful testing will prove which of the above methods works best for you or your particular project!

#### **Use of COMPRESSED DATA = faster computations**

The basic premise is that, because we normally over sample our image data when we digitize it so the pixel size is 3-4 times as fine as the expected resolution, then we can realize significant computational effort merely by performing operations on compressed data. Hence, the program switch, BIN, was added to signal PFTSEARCH to perform its computations on uncompressed data (BIN=1) as before or much faster on compressed data (BIN=2). Recent improvements in memory management have even provided modest program speed-up when PFTSEARCH is run with BIN=1. If BIN=2 (the only other option), then the 3D map is compressed ~8-fold and hence requires a corresponding decrease in memory requirements at run time. Because the "volume" of the compressed 3D map is ~8-fold smaller than that of the original map, 2D projections of the map are computed about 8 times faster. The 'raw' image data are also compressed as they are read in, which speeds up by ~4-fold all manipulations of images and PFT data. The overall improvement in program performance is about 6-8 fold when BIN=2.

## **B. GENERAL DESCRIPTION OF PROGRAM**

PFTSEARCH performs two types of orientation search. They involve coarse, GLOBAL searches (MODE=|1| or |2|). Note that negative values for MODE (-1 or -2) are used to signify that the model projection and PFT data have already been calculated and

stored on disk (using EMPFTRJ) and don't need to be calculated by PFTSEARCH but just read in. In GLOBAL mode operation, the orientation search for each particle image is performed over an entire asymmetric unit, typically in  $1^\circ$  or coarser steps of  $\theta$  and  $\phi$  (see Table 1, later). For MODE=|1| operation (GLOBAL) the particle X,Y origins are unknown or are ignored, and an initial estimate of each particle origin is determined with a real space, cross-correlation procedure. This procedure compares each particle image to a circular average of the sum of all projection views of the model within the asymmetric unit. For MODE=|2| operation, the orientation search makes use of the most recent particle X,Y origin information. In addition, for both MODES of operation, each particle origin is refined after the  $\theta, \phi, \omega$  angles are determined for each particle. Here, a real-space cross-correlation of each particle image is made against the corresponding projected view (at angle  $\theta, \phi, \omega$ ) of the current 3D model.

PFTSEARCH was designed to replace where possible EMICOFV, which uses a common-lines procedure. EMICOFV served for many years as the only means to determine initial orientation parameters for icosahedral particles. **NOTE:** PFTSEARCH does **NOT** replace EMICOFV type routines when there is no *reasonable* model 3D data from which a reference data base can be produced. Thus, EMICOFV (or some variant) may still be needed when initiating a structure determination for a new particle whose structure is unknown or differs substantially from any other previously determined structure and for which no *decent* starting model can be produced. The success of PFTSEARCH in part results because comparisons made between the 'raw' particle images and model projections utilize all the available data, and, in part because the projection data generally have a high signal-to-noise (S/N) ratio. In addition, the PFTSEARCH procedure can be sensitized by selecting just a portion of each image for comparison with the model data. For example, only information within defined radii in real or reciprocal space can be utilized in the calculations.

PFTSEARCH computes two types of data from the 3D model:

1. Projections (**PRJs**) and
2. Polar Fourier transforms (**PFTs**) of the projections.

These are used as reference data to determine  $\theta, \phi, \omega, X, Y$  values for each 'raw' particle image. The high S/N ratios of the model 3D density map and the projections computed from it is what makes the orientation and origin searches work well. However, although this method is quite robust in analyzing images of many unstained, frozen-hydrated specimens, success requires the use of model data that include some elements representative of the correct structure (see e.g. Baker and Cheng [1996] *J. Struct. Biol.* **116:120-130**).

### C. PFTSEARCH PROCEDURE

PFTSEARCH converts each raw image from Cartesian (x,y) to polar (r, $\psi$ ) format by subdividing it into a series of concentric annuli (each usually exactly one pixel wide). Each annulus is sampled

NROT times, with the interpolation performed at a fine enough interval to prevent loss of data resolution at the highest particle radius. Hence, annuli at the lowest radii in the particle image end up way over-sampled. Each annulus is then separately Fourier transformed to give the polar Fourier transform (**PFT**) of the image.

The **PFT** gives the 1D Fourier transform (power spectrum) of the density distribution at each radius of the particle. This provides a particularly useful representation of projection images of spherical objects and allows rapid and sensitive determination of each particle  $\theta, \phi, \omega$  orientation. Program efficiency arises because the particle parameter determination is segregated into a few, discrete stages rather than trying to perform all computations in a single, brute force step.

When  $\text{MODE}=|1|$ , the first step is to estimate the X,Y origin for each particle image by cross-correlating each image with a circular average of the sum of all model projections in the data set (in GLOBAL mode, these cover the entire asymmetric unit). NOTE: the origin of the projection average falls precisely on the central pixel of the NCOL by NROW projection array. Next, in  $\text{MODE } |1|$  or  $|2|$ , the values of  $\theta$  and  $|\phi|$  are determined. Ambiguity in the sign of  $\phi$  corresponds to an ambiguity in knowing the handedness of the particle. Then, the sign of  $\phi$  and the  $\omega$  angle are determined by use of a rotational cross-correlation algorithm which compares each 'raw' image with the best matched projected view of the model and its enantiomer ( $\theta, \phi$  and  $\theta, -\phi$ ) in all possible  $\omega$  orientations. For images of spherical particles this correlation is typically computed only for a restricted range of particle annuli, from ANN\_LO to ANN\_HI. ANN\_HI is generally set to a value close to or just below the outermost radius of the particle. If ANN\_HI is set too large, added background noise from outside the particle reduces the sensitivity of the search procedure. Also, setting ANN\_LO  $\gg$  0 can help in the analysis of images of spherical viruses because the central portions of the projected images of such specimens typically contain a significant contribution from non-icosahedrally ordered material (e.g the genome) and this reduces the effectiveness of identifying the orientation of the icosahedrally-ordered portions of the structure.

The PFT correlation procedure works best when the 'raw' image data are band pass filtered in Fourier space (RES\_LO, RES\_HI) to remove high spatial frequency noise (RES\_HI) and to remove the spherically symmetric, low frequency components (RES\_LO). PFT data, which are 2D, include both real and reciprocal space information. These are usually represented or displayed (ROBEM program) with radius (r) increasing along the ordinate (vertical) axis and Fourier Bessel order (n) increasing along the abscissa (horizontal) axis.

As previously stated, GLOBAL searches are mainly designed to perform a coarse scan of the entire asymmetric unit as a means to 'jump-start' the orientation refinement process. Generally, one cycle each in  $\text{MODE } |1|$  and  $\text{MODE } |2|$  ought to be sufficient to get started so subsequent analysis can be done with the REFINE option ( $\text{MODE}=3$ ). Tips concerning how to use PFTSEARCH (along with EM3DR)

are discussed in greater detail in section **F. Strategy for use of PFTSEARCH and EM3DR.**

**Table 1:** Positional accuracy of orientation search as function of particle size

D(Å )	$\Delta^\circ$				
	2.0°	1.0°	0.5°	0.25°	0.10°
250	8.7	4.4	2.2	1.1	0.43
500	17.5	8.7	4.4	2.2	0.87
1000	34.9	17.5	8.7	4.4	1.75
2000	69.8	34.9	17.5	8.7	3.49

The values in the Table (in Å) represent the distance that points on opposite sides of a spherical object (diameter, D), appear to move in the projection plane when the object is rotated  $\Delta^\circ$ . The following formula is used to compute this distance:

$$\text{Point separation} = \frac{2\pi D\Delta}{360} \text{ \AA}$$

#### D. PROGRAM OUTPUT

The primary output from (also input to) PFTSEARCH is stored in one or more PARAM files. These contain basic information about each scanned IMAGE file (pixel size, voltage, etc.) as well as information about each particle image (ID, THETA, PHI, OMEGA, FFT\_ORIGX, FFT\_ORIGY, MAG\_FAC, and three different correlation coefficients, PFTCC, PRJCC, CMPCC). In addition, various statistics on the correlation coefficients, as a function of radius (PFT.RADS) and as a function of resolution (PFT.RES1 and PFT.RES2) can be generated and stored in output files.

Three correlation coefficients (PFTCC, PRJCC and CMPCC) are listed for each particle in the PARAM file. The first, PFTCC, is a global correlation coefficient computed between the PFTs of each 'raw' image and its corresponding model projection. PFTCC incorporates the effects of any radii cutoffs in real-space and resolution limits in reciprocal space. PRJCC is a global, real-space correlation coefficient computed directly between each 'raw' image (in polar coordinate format) and the corresponding model projection (also in polar format). For PRJCC, the coefficient is only computed for data between radii PFTRAD\_LO and PFTRAD\_HI. This coefficient is computed from polar real space data. The third correlation coefficient, CMPCC, is similar to PRJCC in that the computations are performed on real space data, but CMPCC uses normal Cartesian (x,y) format image and model projection data. **Note:** the values of PRJCC and CMPCC differ, in part (*need to check this*) because low radius data are overweighted in the polar representations of images and model projections. Though opinions are mixed, the CMPCC correlation coefficient *seems* to give the most reliable assessment of particle quality because PFTCC and PRJCC are

more designed to help PFTSEARCH screen for good  $\theta$  and  $\phi$  values (PFTCC) and for good X,Y, and  $\omega$  values (PRJCC).

### **E. BATCH JOB PROGRAM INPUT**

1. MODE, BIN, SYM, DANG, MODE2, ILIST, FILT\_FAC (3I,F,2I,F)
2. PFTRAD\_LO,PFTRAD\_HI,PFTRAD\_STEP,ANN\_LO,ANN\_HI, RADIUS, BFACTOR (3F,3I,F)
3. RES\_LO, RES\_HI, JCUT, SIGCUT, IC, SC, IO, SO (2F,I,5F)
4. MAG\_CEN, MAG\_STEP, MAG\_NUM, MAG\_NORM (2F,2I)
5. 3D model input filename (A) ! DUMMY INPUT IF MODE < 0
6. PRJs input filename (A) ! DUMMY INPUT IF MODE > 0
7. PFTs input filename (A) ! DUMMY INPUT IF MODE > 0
8. empft.rads output filename (A)
9. empft.res1 output filename (A)
10. empft.res2 output filename (A)
11. PARAMETERS filename(s) (A)

## Line-by-line descriptions of program input:

### 1. MODE, BIN, SYM, DANG, MODE2, ILIST, WIENER (3I,F,2I,F)

MODE = 1 Global orientation search mode (no input origins)  
= -1 Global mode (use stored PFTs/PRJs; no input origins)  
= 2 Global mode (use input origins)  
= -2 Global mode (use stored PFTs/PRJs; use input origins)

**CAUTIONARY NOTE:** If you are reading in PFT/PRJ data stored on disk (generated by EMPFTPRJ), then it is **imperative** that you specify the same values for BIN, SYM, DANG, RAD\_LO, RAD\_HI, and RAD\_STEP that you used in EMPFTPRJ. Otherwise, the program is likely to crash or perform other unseemly (and unwanted) calculations.

BIN = 1 Binfactor of 1 (DEFAULT; *i.e.* no data compression)  
= 2 Binfactor of 2 (compresses 3D data eightfold and 2D data fourfold)

The current limit on BIN is 2.

SYM = 1 for no symmetry (like a ribosome)  
= 2-31 for n-fold cyclic symmetry (about z-axis)  
= 532 for 532 icosahedral symmetry (DEFAULT)  
= 222 for 222 dihedral symmetry  
= 32 for 32 dihedral symmetry  
= 422 for 422 dihedral symmetry  
= 52 for 52 dihedral symmetry  
= 622 for 622 dihedral symmetry  
= 72 for 72 dihedral symmetry  
= 822 for 822 dihedral symmetry  
= 92 for 92 dihedral symmetry

SYM specifies the point group symmetry of the particle you are studying. The program searches one-half of the asymmetric unit appropriate for the chosen symmetry. The available symmetries are listed above.

In all search modes the half asymmetric unit (ASU) defines the limits of the search area. Orientations that stray outside the half ASU are folded back to the equivalent view within the ASU. The limits of the half ASU for various symmetries are as follows:

Symmetry	$\theta$ min	$\theta$ max	$\phi$ min	$\phi$ max
1	0.0	90.0	-180.0	180.0
5	0.0	90.0	-36.0	36.0
422	0.0	90.0	0.0	45.0
52	0.0	90.0	0.0	36.0



532	69.09	90.0	0.0	31.71
-----	-------	------	-----	-------

DANG is used to specify the step size in the  $\theta$  and  $\phi$  directions (in degrees: DEFAULT = 1.0). The step size in the  $\theta$  direction remains constant (= DANG) but it varies in the  $\phi$  direction from a smallest value (= DANG) when  $\theta=90^\circ$  (at the 'equator') and *increases* thereafter for progressively smaller values of  $\theta$ . Varying the  $\phi$  step size assures uniform sampling of the ASU in regions where  $\theta$  is  $<90^\circ$ . The step size in the  $\phi$  direction is given by the formula:

$$\text{DANG}/\sin\theta$$

If the  $\phi$  step size was not varied as given in the above formula, the grid sampling would be much too fine near the 'poles'. When  $\theta$  approaches  $0^\circ$  or  $180^\circ$  (N and S poles), the program adjusts the  $\phi$  step size to maintain even sampling. For example, at  $\theta = 30^\circ$ , the  $\phi$  step size will be twice DANG (*i.e.*  $2.0^\circ$  if DANG is  $1.0^\circ$ , since  $\{1.0/\sin(30^\circ)\}=2.0$ ).

The value you choose for DANG has a **tremendous** impact on PFTSEARCH run time. As Table 1 showed (**Section C.**), it makes no sense to set this parameter very small, especially at the beginning of data analysis when GLOBAL orientation searches are performed. Table 2 demonstrates how the number of program calculations significantly increases as DANG is decreased. Also, Table 2 clearly demonstrates how the "problem" becomes even more severe in the case of particles with lower symmetry.

**Table 2:** Number of views per ASU as a function of particle symmetry and orientation search angle increment ( $\Delta^\circ$ )

$\Delta^\circ$	Point Group Symmetry			
	532	52	5	1
3.0 0	48	266	519	2,353
2.0 0	91	573	1,127	5,253
1.0 0	370	2,173	4,309	20,809
0.5 0	1,430	8,471	16,869	82,869
0.2 5	5,606	33,439	66,733	330,751
0.1 0	34,327	207,358	414,353	2,064,448

MODE2 determines if and how the CTF is used. **Note:** It only makes sense to use this option **IF** the 3D model is a CTF-corrected map. If the map is an uncorrected one, then the use of MODE2 = 1 or 2 will force the program to compare CTF-modified and CTF-unmodified data, which is undesirable and may lead to undetermined errors in refinement. The values of MODE2 may be set as follows:

- MODE2 = 0 No CTF modifications to the data are made.
- = 1 Projections of the 3D model are multiplied by the CTF, and these are compared to the unmodified, raw images.
- = 2 The raw images are multiplied by the CTF, and these are compared to unmodified projections of the 3D model.
- = 3 The raw images are multiplied by the INVERSE CTF, and these are compared to unmodified projections of the 3D model.

ILIST signals PFTSEARCH to generate various forms of output:

- ILIST = 0 Minimal output (refined orientations and origins)
- = 1 Correlation coefficients are computed as functions of radius and resolution and results are stored in output files PFTSEARCH.RADS, PFTSEARCH.RES1, PFTSEARCH.RES2.
- = 2 Same as ILIST=1, but also gives correlation coefficients for each particle view

## **2. PFTRAD\_LO, PFTRAD\_HI, PFTRAD\_STEP, ANN\_LO, ANN\_HI, RADIUS, BFACTOR (3F,3I,F)**

These variables (real space pixel units) define the annular portion of the projected data (PFTRAD\_LO to PFTRAD\_HI) to be used in the PFT calculations. PFTRAD\_STEP sets the radial step interval and hence determines the number of annuli in each PFT ( $NANNULI = \lceil [PFTRAD\_HI - PFTRAD\_LO] / PFTRAD\_STEP \rceil + 1$ ). PFTRAD\_LO is normally left = 1.0 and PFTRAD\_HI is usually set to a value just larger than the particle boundary (but usually much smaller than  $NCOL/2$  if you boxed the original particle images conservatively). Use PFTRAD\_LO > 1.0 if you suspect that the projected data at higher radii (*i.e.* projection only of outer capsid features) will give a more sensitive measure of the orientation parameters.

PFTRAD\_LO can't be set lower than 1.0 because the center of the projected view doesn't change with orientation and hence gives no useful information for the PFT calculations. The default for PFTRAD\_HI is  $(NCOL+1)/2$ , but note that this will generally be too large especially if the particle boxing was performed conservatively.

PFTRAD\_STEP is normally left = 1.0. Using a larger value will reduce the number of computations (smaller NANNULI) but use cautiously (especially if BIN=2) so you don't sample the data too coarsely. If PFTRAD\_STEP is < 1.0, you may make needless calculations (NANNULI too large), especially if BIN=1 and the images were digitized at a pixel resolution at least twice as fine as the expected resolution limit. NOTE: This variable is slated for removal in the future.

ANN\_LO, ANN\_HI define the range of annuli from the real space polar coordinate image and projection data that are included in the calculations. The best way to determine optimum values for these

parameters is to run PFTSEARCH with the initial DEFAULT values (0,NANNULI-1), and check the output file PFTSEARCH.RADS to see the correlation coefficients a function of radius. The average correlation coefficient typically undergoes a large drop near the outer edge of the particle (ANN\_HI). Also, the correlations are generally lower at low radii (ANN\_LO) corresponding to the 'core' part of the structure. Thus, ANN\_LO and ANN\_HI may require some fine tuning to optimize the orientation and origin search procedure.

The chosen bin value (Line #1) does not change the entered value of PFTRAD\_LO, PFTRAD\_HI, ANN\_LO and ANN\_HI. Conversions are made automatically by the program.

**RADIUS and BFACTOR not used by PFTSEARCH, they are placeholders for compatibility with other programs such as OOR.**

### **3. RES\_LO, RES\_HI, JCUT, SIG, IC, SC, IO, SO (2F,I,5F)**

RES\_LO and RES\_HI define the lower and upper resolution limits of the data to be included in the calculations. These program input variables should be specified in the same units as PIXSIZE, where PIXSIZE is the size of each pixel in the digitized images and is specified for each image in the corresponding PARAM file. PIXSIZE may be defined in any units you choose (Å, nanometers, pixels, yards, etc.) but you MUST be consistent and use the same units to define PIXSIZE, RES\_LO and RES\_HI. A value of RES\_HI smaller than 2\*PIXSIZE (the DEFAULT) is disallowed (this would exceed the Nyquist limit which is two-pixel resolution). You should treat this value for RES\_HI as an absolute lower limit, which, if used, is likely to be an unrealistic value for real (*i.e.* noisy) data. Hence, use careful judgment in setting the value of RES\_HI.

The DEFAULT for RES\_LO [IDIM1\*PIXSIZE] is probably unrealistically large. Again, PFT works best if some of the low resolution data are ignored. This is because, for example, for particles with icosahedral symmetry the very low resolution data only carry information about spherical symmetry and little if any information about icosahedral symmetry. General rule of thumb: set RES\_LO to a value no smaller than one-fifth the size of the particle diameter. For example, if the particle diameter is 500Å, then a value for RES\_LO of 100Å might be an appropriate starting point for PFTSEARCH. If you prefer to specify PIXSIZE in pixel units (instead of Å or another unit), and the pixel size of the digitized image corresponds, for example, to 3.86Å units, then RES\_LO would be 25.9 (=500/(5\*3.86)).

**REMINDER:** Use good judgment in setting the above program variables!!! The success or failure of PFTSEARCH may very well depend on your ability or lack thereof to set these values.

JCUT specifies the minimum rotational Bessel order (Jn) to include in the calculations. The DEFAULT (=1) omits the Jo term, which is recommended in analyzing icosahedral particles because

this removes the spherically symmetric image components and boosts the sensitivity of determining icosahedral orientations. To include Jo, set JCUT to any negative value. To cut out higher orders, JCUT is set to a value greater than 1 (NOTE: to my knowledge, no one has ever carefully tested the effects of doing this).

SIG allows the program to filter the PFT data on the basis of the variance of the PFT data. In theory, this option should greatly sensitize PFTSEARCH. However **BE FOREWARNED**: this option is still **UNTESTED** so it may not and probably does not work!!! SIG specifies the threshold for the variance mask. With SIG=0.0 (DEFAULT), the masking option is disabled.

IC, SC, IO and SO not used by PFTSEARCH, they are placeholders for compatibility with other programs such as OOR.

#### 4. MAG\_CEN, MAG\_STEP, MAG\_NUM, MAG\_NORM (2F,2I)

These program input parameters are used to direct the magnification factor refinement and CMP correlation coefficient (CMP\_CC) calculations.

MAG\_CEN specifies the midpoint of the magnification scale factor search. The program tests for magnification factors that are  $(MAG\_NUM-1)/2$  steps above and below MAG\_CEN. Set MAG\_CEN = 0.0 or a negative number to force the program to use the MAG factor for each image as stored in the PARAM input data file. Hence, when any positive value of MAG\_CEN (e.g. 1.0) is chosen, the MAG search "window" for all images will be over the same range. **\*\*\* NEED TO CHECK THIS OUT \*\*\***

MAG\_STEP defines the grid size of the magnification search. A value of MAG\_STEP = 0.005, for example, corresponds to 0.5% increments.

MAG\_NUM establishes the extent of the magnification search "window". This should be an odd integer > 0. For example, with MAG\_CEN = 1.0, MAG\_STEP = 0.005, and MAG\_NUM = 11, the search window will encompass magnification factors ranging between 0.975 and 1.025. Entering '1' will force MAG to be either MAG\_CEN or the value read from the PARAM file (which occurs whenever MAG\_CEN is set ≤ 0.0).

MAG\_NORM is a switch used to normalize the MAG scale factors so that the average MAG is 1.0. This occurs only when MAG\_NORM is set equal to 1, otherwise, the MAG values determined by the program are output without being normalized.

#### 5. 3D model input filename (A FORMAT)

Enter the name of the file that contains the 3D model from which new PRJs and PFTs are to be calculated. **Note**: PFTSEARCH currently only works if the dimension (NCOL\_MAP = NROW\_MAP = NSEC\_MAP) of the 3D model **exactly** matches the image dimension.

## 6. PRJ input filename (A FORMAT)

Enter the name of the file that contains the model projection data (DEFAULT = PFT.PRJS). This data is only read in if MODE = -1 or MODE = -2, otherwise the program calculates what it needs at run time.

## 7. PFT input filename (A FORMAT)

Enter the name of the file that contains the model PFT data (DEFAULT = PFT.PFTS). This data is only read in if MODE = -1 or MODE = -2, otherwise the program calculates what it needs at run time.

## 8. empft.rads filename (A format)

Enter the name of the file that will contain the radial distribution of correlation coefficients.

## 9. empft.res1 filename (A format)

Enter the name of the file that will contain the distribution of cross correlation coefficients PMAP versus POLAR\_PRJ.

## 10. empft.res2 filename (A format)

Enter the name of the file that will contain the distribution of cross correlation coefficients filtered PMAP versus POLAR\_PRJ.

## 11. PARAM input filename(s) (A FORMAT)

Enter the names of up to 999 PARAM files, with each filename on a new input line. The format of each PARAM file is:

LINE INPUT

1 IMAGE filename (A)

This is the name of the IMAGE file which contains byte-packed raw image data stored in \*.PIF format.

2 PIXSIZE, UNITS, VOLTS, AMP\_FAC, DELF\_MAJ, DELF\_MIN, ANG\_MAJ, Cs (F,I,6F)

PIXSIZE = pixel size for data in IMAGE file. May be specified in any units (e.g. Å, nanometers, pixels, etc.) as long as you make sure to specify the correct value for UNITS.

UNITS = specifies the units assigned to PIXSIZE. UNITS is defined as follows:

= 0 assumes PIXSIZE is given in dimensionless pixels

= 1 assumes PIXSIZE is given in Ångstroms

- = 2 assumes PIXSIZE is given in nanometers
- VOLTS = microscope voltage (in volts) for the image data in the IMAGE file.
- AMP\_FAC = amplitude factor. For cryoEM data, a DEFAULT value of 0.07 is reasonable and anything above 0.15 might be suspicious.
- FOCUS\_MAJ = defocus value ( $\mu\text{m}$ ) along the major axis of astigmatism. **Note:** Positive values are used to designate underfocus.
- FOCUS\_MIN = defocus value ( $\mu\text{m}$ ) along the minor axis of astigmatism.
- ANG\_MAJ = angle between major axis and X-direction in FFT, measured positive in a counter-clockwise direction (X-axis = 0.0)
- Cs = spherical aberration coefficient (mm) for the microscope used to record the data in the IMAGE file.
- 3 ID, THETA, PHI, OMEGA, FFT\_ORIGX, FFT\_ORIGY, MAG\_FAC, PFTCC, PRJCC, CMPCC (I,9F)
- ID = specifies the particle # in the \*.PIF format IMAGE file.
- THETA, PHI, OMEGA = orientation of particle #ID
- FFT\_ORIGX, FFT\_ORIGY = center of particle #ID (FFT coordinates)
- MAG\_FAC = scale of particle #ID relative to a standard (model)
- PFTCC, PRJCC, CMPCC = three correlation coefficients
- N+2 Same as #3 for as many particles (N) as needed.

Some parameters may equal 0.0 or may be left blank depending on the stage of the analysis. For example, when PFTSEARCH is first run (MODE=1), the PARAM file will simply consist of the first two lines as outlined above. There will be **ABSOLUTELY NO** lines of particle data (or they will be ignored in MODE=1) because PFTSEARCH performs a global search of **ALL** particle images in the PIF format IMAGE files specified in line #1 of each PARAM file. The first run of PFTSEARCH creates a PARAM file with the first set of ID, THETA, PHI, OMEGA, X, Y, MAG\_FAC, PFTCC, PRJCC, and CMPCC values.

PFTSEARCH outputs a new PARAM file for each one used as input. The new PARAM data files are named with a specific convention: a "\_00#" is tagged after each input PARAM filename. For example, suppose you had but one input PARAM file named "MYDATA.DAT\_000". Then, the name of the output PARAM file would be "MYDATA.DAT\_001". Normally you would use *this* file as the input file for the next run

of PFTSEARCH, in which case the subsequent PARAM file would automatically be named "MYDATA.DAT\_002", and so forth. The importance of understanding the way this works is **crucial!** If, by 'accident', you forget to update the correct PARAM filename in your BATCH COMMAND PROCEDURE file prior to the next cycle of PFTSEARCH (e.g. leaving the name "MYDATA.DAT" in the command file), the output will end up in a file named "MYDATA.DAT\_001" **INSTEAD** of "MYDATA.DAT\_002".

**FINAL NOTE:** if you use several PARAM files as input to PFTSEARCH, then use some rational naming system so the output files will be clearly distinct from the input files. Here's one example in which hypothetical PARAM files for four different micrographs are used:

<u>PFTSEARCH</u>		
<u>Cycle#</u>	<u>Input PARAM filenames</u>	<u>Output PARAM filenames</u>
1	HSV_1856.DAT_000	HSV_1856.DAT_001
	HSV_1858.DAT_000	HSV_1858.DAT_001
	HSV_1900.DAT_000	HSV_1900.DAT_001
	HSV_1989.DAT_000	HSV_1989.DAT_001
3	HSV_1856.DAT_001	HSV_1856.DAT_002
	HSV_1858.DAT_001	HSV_1858.DAT_002
	HSV_1900.DAT_001	HSV_1900.DAT_002
	HSV_1989.DAT_001	HSV_1989.DAT_002
2	HSV_1856.DAT_002	HSV_1856.DAT_003
	HSV_1858.DAT_002	HSV_1858.DAT_003
	HSV_1900.DAT_002	HSV_1900.DAT_003
	HSV_1989.DAT_002	HSV_1989.DAT_003

```

!*****
!* Example PFTSEARCH BATCH JOB run in GLOBAL MODE with      *
!* PFT data read in from disk file (MODE = -1) and          *
!* using compression factor of 2 (BIN = 2) and              *
!* normalizing the MAG scale factors (MAG_NORM = 1)         *
!* and using data from three PARAM files.                   *
!*****
-1, 2, 532, 1.0, 0, 0
1.0, 46.0, 1.0, 0, 46, 0, 200.0
200., 40., 1, 0.0, 0.7, 4.0, 0.9, 3.0
0.0, 0.005, 11, 0
HSV.PIFMAP
HSV.PRJS
HSV.PFTS
Empft.rads
Empft.res1
Empft.res2
HSV_1856.DAT_000

```

```
HSV_1858.DAT_000  
HSV_1989.DAT_000
```

```
Point to appropriate directory:  
% cd ~tsb/v/hsv/test
```

```
Run pftsearch with the file created above ("pftsearch.bch") and  
create log file "pftsearch.log":  
% pftsearch < pftsearch.bch > pftsearch.log &
```



## F. STRATEGY FOR USE OF PFTSEARCH, OOR AND EM3DR

The complexity of computing a 3D reconstruction from a set of cryoEM images means that there is no simple *cookbook* strategy for accomplishing this task. Each set of image data presents its own challenges. Users must be diligent and devise an appropriate strategy based on several criteria such as size of particle, magnification of images(s), digitization step size, resolution range(s) to use, type of specimen (is it a brand spanking new structure never seen before or one that may be new but closely related to a virus whose structure is already known?), etc. The day is yet to come when programs are able to take as input a set of images and spit out a *correct* 3D structure! However, ample experience indicates that the following, general strategy may prove useful in many applications:

1. Obtain a 3D density map to use as a starting model.
2. Run PFTSEARCH (GLOBAL MODE=1) with BIN=2 and MAG\_CEN=1.0. Use fairly coarse search interval (almost certainly *no smaller* than  $1^\circ$  - see Table 1) to obtain the first set of  $\theta, \phi, \omega, X, Y$  parameters.
3. Run EM3DR (BIN=2) to compute a new 3D map from images with the set of  $\theta, \phi, \omega, X, Y$  parameters obtained in Step #2.
4. Run PFTSEARCH (GLOBAL MODE=2) with BIN=2, using the same search interval as in step #2 to obtain new  $\theta, \phi, \omega, X, Y$  parameters.
5. Compare latest  $\theta, \phi, \omega, X, Y$  parameters to those found in Step #2. Examine the three different correlation coefficients (see below) to detect 'bad' particle images and to compare with the coefficients found in Step #2 (they *should*, on average, be better!).
6. Run EM3DR (BIN=2) to compute a 2nd 3D map from images with the  $\theta, \phi, \omega, X, Y$  parameters obtained in Step #4.
7. Run OOR (REFINE MODE=3), using the same search interval as in Step #2.
8. Repeat steps #6 and 7 for several cycles. With each new cycle you need to decide whether or not to change program parameters (it is probably best to stick to ONE change at most per cycle). You may consider changes such as:
  - a) Reduce search step size in half. For example, if the step size was  $1^\circ$  in the previous cycle, reduce it to  $0.5^\circ$  for the next cycle.
  - b) Increase the resolution limits of the PFTSEARCH search and/or EM3DR output (use RES\_LO and RES\_HI variables, see below).
  - c) Filter out 'bad' particle images from EM3DR.
9. Repeat step #8 for one or two cycles with BIN=1. This is done mainly to make sure everything is OK. Of course, cycles with BIN=1 take considerably longer to execute than previous ones because uncompressed data are used in the calculations. If you

notice any drastic changes in correlation coefficients as listed by PFTSEARCH or the  $\theta, \phi, \omega, X, Y$  parameters suddenly change, then you may have proceeded too quickly during previous cycles and may need to back up a few steps before continuing on.

## G. DISCLAIMER

To date (5/98) this program has worked quite well in the analysis of a wide variety of specimen images. The routines seem to work as long as the starting model as well as subsequent reconstruction maps are on track. If your starting model is seriously flawed, don't expect miracles! Also, for particle image data that exhibits very weak enantiomorphic features, the model must include some correct enantiomorphic character or the refinement will lead to a 3D map that exhibits mirror line symmetry about the equatorial line. In our experience this has not been a problem for icosahedral particles with handed surface lattices (e.g. T=7 papovaviruses) because the arrangement of morphological units is clearly enantiomorphic even at very low resolution (>50Å). In many instances, enantiomorphic features do not become apparent until much higher resolutions (<30-40Å), and therefore proper refinement of data cannot proceed until the model incorporates information at the higher resolutions. In tricky situations, it is still advisable to use programs like EMICOGRA with small data sets (5 particles or less) to try and make sure that the particles are oriented with respect to a consistent hand. A crude 3D reconstruction computed from such a limited data set, although noisy, may give a much better model for further refinement with PFTSEARCH.

## H. CREDITS

The original code for PFTSEARCH was developed by T. Baker in ~1988 and tested in a class project by J. Tesmer in ~1990. R. H. Cheng performed more extensive and rigorous tests which led to a full scale working version of the program in 1993. A description of the program as used during the period from about 1993 up through 1997 appears in:

Baker, T. S. and R. H. Cheng (1996) A model-based approach for determining orientations of biological macromolecules imaged by cryoelectron microscopy. *J. Struct. Biol.* **116:120-130**.

Some preliminary discussion of the routine also appears in:

Cheng, R. H., V. S. Reddy, N. H. Olson, A. J. Fisher, T. S. Baker, and J. E. Johnson (1994) Functional implications of quasi-equivalence in a T=3 icosahedral animal virus established by cryoelectron microscopy and X-ray crystallography. *Structure* **2:271-282**.

A number of people have and continue to make valuable additions/corrections/suggestions to PFTSEARCH. These include (among others and in alphabetical order), Robert Ashmore, Steve Walker, Wei Zhang at Purdue, David Belnap and James Conway (NIH), and Stephen Fuller (EMBL).



## K. FLOW CHART FOR PFTSEARCH PROGRAM

\*\* means "see below"; # means end of program; ! signifies no subroutine calls

```
*****
*   MAIN   *
* (PFTSEARCH.FOR) *
*****
*
*   | - GET_NVIEW_MOD - GET_T1T2P1 !
*   | - PIF_OPEN !
*   | - PIF_READ_GH - differentEndian !
*- INFO - | - PIF_READ_DH -----| - differentEndian !
*   |                                     | - convertBackFloat !
*   |
*   | - PFILE_INFO - | - PIF_CLOSE !
*   |                 | - PIF_OPEN !
*   |                 | - PIF_READ_GH - differentEndian !
*   | - PIRADDEG !
*   | - FFT_SETDIM_DEF_SAME !
*
*- GETMEM0 - MALLOC !
*
*- FILL_PARAMS !
*
*   | - GETMEM1 - MALLOC !
*   | - GET_T1T2P1 !
*   | - CALC_MOD_TPS !
*   | - GETMEM2 - MALLOC !
*   | - PIF_READ_WMAPI2 - PIF_READ_MAPI2**
*   | - PIF_CLOSE !
*   | - STOP_WATCH_START !
*- GLOBAL - | - GET_PRJSPFTS_G**
*   | - STOP_WATCH_STOP !
*   | - FREEMEM2 - FREE !
*   | - GETMEM3 - MALLOC !
*   | - PFTCC_FILL_G - | - PIRADDEG !
*   |                 | - PIF_READ_MAPI4**
*   | - GLOBAL_CC**
*   | - FREEMEM1 - FREE !
*   | - FREEMEM3 - FREE !
*
*- PARAM_SRT (Lixun's routine)
*
*   | - STOP_WATCH_START !
*   | - GETMEM_REFINE - MALLOC !
*   | - PIF_READ_WMAPI2 - PIF_READ_MAPI2**
*   | - PIF_CLOSE !
*   | - CALC_TP_MASK - PIF_WRITE_DEBUGI4 - PIF_WRITE_BYTE_IMAGE !
*   | - PIF_WRITE_DEBUGI4 - PIF_WRITE_BYTE_IMAGE !
*   | - GETMEM_REFINE2 - MALLOC !
*
*- WRITE_PARAMS - INCREMENT_FNAME !
*- PIF_CLOSE !
*- STOP_WATCH_STOP !
*- CALC_AVGS !
*- FREEMEM0 - FREE !
#
```

```
=====  
=  
|  
| | | - PIRADDEG !  
| | | - MAP_CLEAR !  
| | - MAP_PRJ - | - MAP_PRJ_XZ --- MAP_CLEAR !  
| | | - MAP_PRJ_AXIS - MAP_CLEAR !  
| | | - MAP_PRJ_ALL -- MAP_CLEAR !  
|  
| | - SAVE_PRJS - MAP_SYM**  
- GET_PRJSPFTS_G - | - MAP_POLAR - | - MAP_CLEAR !  
| | | - MAP$POLAR - MAP_POLAR_GRID !  
|  
| - PMAP_FFT - FOURT - L6TO9 !  
=====  
=
```

```
|- PIRADDEG !  
|- GETMEM4 - MALLOC !  
  
| | | - MAP_CLEAR !  
| | | - PIF_READ_MAPI4**  
- PRJAVG_FFT - | - MAP_FFT_FILL - FFT_CLEAR !  
| | | - FFT_2D - FOURT - L6TOL9 !  
  
- PIF_OPEN !  
- PIF_READ_GH - differentEndian !  
- PIF_CLOSE !  
- PIF_READ_DH - | - differentEndian !  
| | | - convertBackFloat !  
  
- PIF_READ_IMG11 - PIF_READ_BYTESHORT_IMAGE !  
- IMG_MAP !  
- MAP_FFT_FILL - FFT_CLEAR !  
- FFT_2D - FOURT - L6TOL9 !  
- CTF_MULTIPLY - CTF_ASTIG !  
- FFT_2D_BACK - FOURT - L6TOL9 !  
- FFT_MAP_FILL !  
- MAP_FFT_FILL - FFT_CLEAR !  
- FFT_2D - FOURT - L6TOL9 !  
- GET_XY**  
- INTERP2D (NIH C routine)  
- MAP_FFT_FILL - FFT_CLEAR !  
- FFT_2D - FOURT - L6TOL9 !  
- INTERP2D (NIH C routine)  
- MAP_FFT_FILL - FFT_CLEAR !  
- FFT_2D - FOURT - L6TOL9 !  
  
| | | - MAP_POLAR - | - MAP_CLEAR !  
| | | | - MAP$POLAR - MAP_POLAR_GRID !  
|- GLOBAL_CC - | - GET_TPO_G - | - PMAP_FFT - FOURT - L6TO9 !  
| | | | - AVG_PFTIMG !  
| | | | - GET_THEPHI_G !  
| | | | - GET_BESTPRJ - PIF_READ_MAPI4**  
| | | | - GET_PHIOMEGA !  
|- MAP_XFLIP !  
|- COPY_R4 !  
|- INTERP2D (NIH C routine)
```

```

- MAP_FFT_FILL - FFT_CLEAR !
- FFT_2D - FOURT - L6TOL9 !
- GET_XY**
    | - IMG_FFT_FILL - FFT_CLEAR !
    | - MAP_FFT_FILL - FFT_CLEAR !
    | - FFT_2D - FOURT - L6TOL9 !
    | - FFT_LOP !
- GET_BESTMAG - | - FFT_HIP !
    | - FFT_2D_BACK - FOURT - L6TOL9 !
    | - FFT_IMG_FILL !
    | - FFT_MAP_FILL !
    | - INTERP2D (NIH C routine)
- LIST_CCS**
    | - TO_NFOLD !
- TO_ASYM_UNIT - | - RETURN_TO_UNIT_TRIANGLE**
    | - TO_N22 !

- FREEMEM4 - FREE !

```

=

```

- LIST_CCS -
    | - COPY_I2 !
    | - IMG_MASK !
    | - IMG_FFT_FILL - FFT_CLEAR !
    | - FFT_2D - FOURT - L6TOL9 !
    | - COPY_R4 !
    | - MAP_2DMASK !
    | - MAP_FFT_FILL - FFT_CLEAR !
    | - PFTCC_RES !
    | - FFT_FLT - | - FFT_HIP !
    |               | - FFT_LOP !
    |
    | - FFT_2DBT - | - COPY_R4 !
    |               | - FOURT - L6TOL9 !
    |
    | - IMG_MAP!
    |
    | - MAP_POLAR - | - MAP_CLEAR !
    |               | - MAP$POLAR - MAP_POLAR_GRID !
    |
    | - PFTCC_RAD !

```

=

```

- GET_XY - | - CCF_FFT - CCF - | - MAP_CLEAR !
    |               | - FOURT - L6TOL9 !
    |               | - MAP_STATS - MAP$STATS !
    |
    | - PFTCC_PEAK - | - MAP_PEAK !
    |               | - TRUS_POLYS - TRUS_SOLVE - TRUS_MATINV !

```

```

=====
=
|
| - MAP_SYM_CAVG - MAP_STATS - MAP$STATS !
| - MAP_SYM_RAVG - | - MAP_SYM_GRID !
- MAP_SYM - | - COPY_R4 | - MAP_STATS - MAP$STATS !
| - MAP_STATS - MAP$STATS !
=====
=
|
| - POLAR_TO_COSINES - NORM !
- RETURN_TO_UNIT_TRIANGLE - | - EQUIVALENT_VIEW - | - CROWTHER_TO_MATRIX !
| - MATMUL !
| - MATRIX_TO_CROWTHER !
=====
=
|
| - PIF_READ_MAPI4 - PIF_READ_MAP_INT_IMAGE - differentEndian !
=====
=
|
| - PIF_READ_MAPI2 - PIF_READ_MAP_SHORT_IMAGE - differentEndian !
=====
=
|
| - PFTCC_R (***) NOT WRITTEN YET (***)
|
=====
=

```