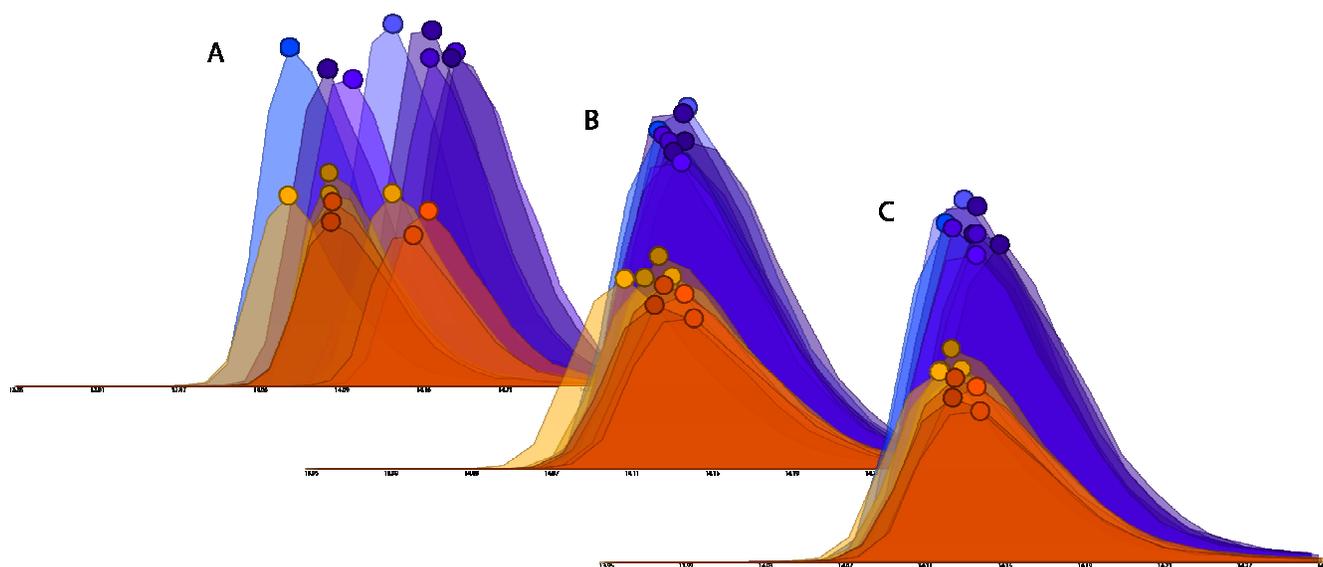


## Pteroylglutamic acid(Folate) m/z: 440.1300-440.1348



**Supplementary Figure 1: Two iterations of alignment.** Shown are sixteen EICs in the m/z and retention time range of folate from a viral infection time course. (A) Before alignment. (B) First iteration of alignment. The alignment was based on 2000 high intensity peaks (>5000 ion counts) found in all samples. (C) Second iteration of alignment. Complete list of alignment parameters used is listed below.

### Peak Selection Criteria (peaks that are used to create the alignment):

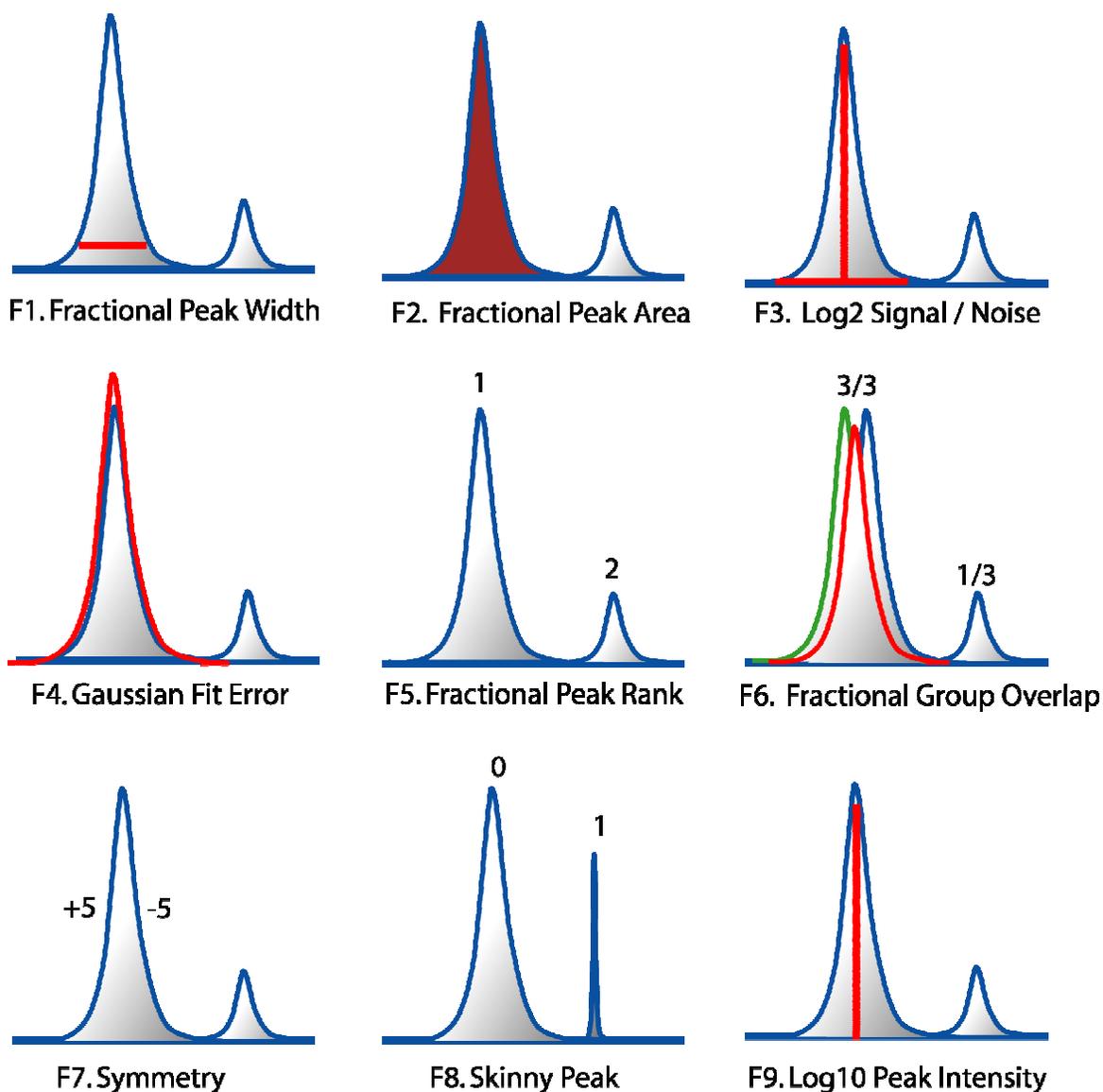
1. Minimum Peak Intensity is 5000 ions
2. Minimum Peak S/N Ratio is 5
3. Minimum Peak Width is 5 scans

### Peak Grouping (peak groups that are used to create the alignment):

1. Peak grouping retention time window is 20 scans (peaks that are within  $\pm 20$  scans of the median peak center are grouped together)
2. Group must contain at least 2 good peaks.
3. Maximum number of groups used in the alignment is 2000 (groups are ranked by intensity)

### Alignment Algorithm:

- Maximum number of alignment iterations is 10
- Fitting using polynomial of degree 3



**Supplementary Figure 2: Peak features used for quality scoring.** F1: Fractional Peak Width = peak width at baseline normalized by total number of scans in EIC with nonzero intensity; F2: Fractional Peak Area = peak area normalized by total area of EIC; F3: Log2 Signal to Noise (S/N) = Standard S/N ratio, normalized by maximum S/N ratio found for any compound and log2 transformed; F4: Gaussian Fit Error = Root mean square difference between observed shape and best fitting Gaussian curve; F5: Fractional Peak Rank = peak rank upon sorting peaks within EIC by area, normalized by number of peaks in the EIC; F6: Fractional Group Overlap = number of overlapping peaks from other samples normalized by the total number of samples; F7: Symmetry = count of scans increasing in intensity plus of scans decreasing in intensity scaled by peak width as per Formula 1 below; F8: Skinny Peak = True (1) if peak width is less than 3 scans, else false (0); F9: Log10 Peak Intensity = log10 of peak height.

Peak Symmetry score =  $(nSteps - \text{abs}(nPos - nNeg)) / nSteps * \log_2(nSteps)$  where:

nPos : count of scans increasing in intensity

nNeg: count of scans decreasing in intensity

nSteps: nPos + nNeg;

Example: Symmetric peak, with +5 and -5: score = 3.3

Asymmetric peak with +1 and -5 : score = 0.8

Mass slice detection, EIC extraction, Peak detection, Peak scoring, and Peak grouping Benchmark  
 Computer: ( Intel 6600 2.4ghz, 2 cores with 8 gb ram)

**Peak Detection Settings**

<b>Mass Slice Detection</b>			
Mass Domain Resolution	20 ppm		
Time Domain Resolution	10 scans		
<b>EIC Processing</b>			
Smoothing Window	10 scans		
Peak Grouping Window	0.5 min		
			12.5
			11.875
<b>Peak Scoring Parameters</b>			
Minimum Number of "good" peaks in group	2 peaks		11.875
Minimum Signal/Baseline Ratio	2		
Minimum Peak Width	3 scans		
Minimum Signal/Blank Ratio	2		
Minimum Peak Intensity	5000 ions		

**Results**

	#Files	Memory Usage	Run Time	#Mass Slices	#EICs	#Peaks	#Groups
<b>Using One Core</b>							
Dataset: 16 centroided mzXML files, ~12 Mb each	16	300 MB	21.09 sec	5628	90,000	137,295	5680
Dataset: 32 centroided mzXML files, ~12 Mb each	32	480 MB	46.00 sec	7132	228,000	346,296	6561
Dataset: 48 centroided mzXML files, ~12 Mb each	48	670 MB	64.0 sec	7086	340,000	523,878	6603
<b>Using Two Core</b>							
Dataset: 16 centroided mzXML files, ~12 Mb each	16	300 MB	11.06 sec	5628	90,000	137,295	5680
Dataset: 32 centroided mzXML files, ~12 Mb each	32	480 MB	28.08 sec	7132	228,000	346,296	6561
Dataset: 48 centroided mzXML files, ~12 Mb each	48	670 MB	38.05 sec	7086	340,000	523,878	6603

**Supplementary Table 1:** Benchmark of Maven memory and CPU usage. The memory and speed of the program was benchmarked on three sets of files (16 samples, 32 samples, and 48 samples) from viral infection dataset. All files were loaded into memory and peak detection algorithms were executed using settings shown in the table above. Baseline memory consumption of the program (physical memory used without any files loaded) was ~ 120 MB, with additional 12Mb of usage for each loaded file. The performance was evaluated using the same computer operating on a single CPU and operating on two CPU cores. We find that overall performance scales roughly linearly with number of samples and number of CPUs used.