

# Compression of IASI Data and Representation of FRTM in EOF Space

Peter Schlüssel  
EUMETSAT  
Am Kavalleriesand 31  
64295 Darmstadt  
Germany

Anthony C. L. Lee  
Met Office  
FitzRoy Road  
Exeter, Devon EX1 3PB  
United Kingdom

# Outline

- Current baseline
- Possible approach for IASI data compression
- Apodisation
- Cloud detection using EOF scores
- RTM formulation in EOF space

# Current Baseline

- The data of the Infrared Atmospheric Sounding Interferometer (IASI), resolving the spectrum between 645 and 2760  $\text{cm}^{-1}$  at 0.5  $\text{cm}^{-1}$ , are sampled at 0.25  $\text{cm}^{-1}$  intervals, thus representing 8461 spectral samples
- The disseminated Level 1c spectra are apodised to standardised spectral response function with a Gaussian shape and a half width of 0.5  $\text{cm}^{-1}$ 
  - Same ISRF for all channels
  - Avoids negative radiances which occur in self-apodised spectra
- The day-1 approach is to disseminate IASI Level 1c data as 8641 spectral samples quantised to 16 bit/sample

# Current Baseline (cont.)

- **Advantages of Day-1 format**

- Users can handle the IASI spectra in the same way as data from traditional channel radiometers by picking individual samples (“channels”) to support their particular application
- Peculiarities of the interferometric measurements, like negative radiances, are hidden

- **Disadvantages of Day-1 format**

- The data volume is bulky: 2 Mbit/s
- Quantisation in 16-bit samples will slightly degrade the spectra
- Exploitation of full information contained in IASI data to support NWP forecast is prohibited due to huge number of spectral samples
- Apodisation of spectra introduces non-diagonal error covariance, which complicates (and may inhibit) use of adjacent and nearby channels in data assimilation

# Data Volume Reduction

- **Data Thinning** can be done spectrally by selecting only a sub-set of channels for transmission to users, or spatially, by communicating only a horizontally sub-sampled set of soundings. The latter bears the danger of losing meteorologically interesting situations
- **Principal Component Analysis** allows the projection of the spectra on to a pre-defined set of eigenvectors. Of the corresponding scores a set of carefully truncated sub-sets can be communicated, from which most, though not all, of the spectrum can be re-constructed
- **Compression** of the data by means of run-length and Huffman encoding, adapted to the IASI data characteristics, is possible after carefully controlled quantisation.

# IASI Data Representation

- Before the data are communicated to the users it has to be represented at controlled quantisation, within a suitable dynamic range
- **Dynamic Range:** The Day-1 processing limits the data to the 180 K to 315 K range; for a compressed format (see below) no restriction is necessary
- **Quantisation:** This must be related to the instrument noise (typically described as  $NE\Delta T$ )
  - Radiance quantisation adds white, or uniform, radiance noise-power spectral density
  - RMS amplitude of added quantisation noise depends on quantisation step size
  - Steps in fractions of 0.5, 1.0, 2.0, 4.0  $NE\Delta T$  lead to  $NE\Delta T$  increase of 1.0, 4.1, 16, 53%, i.e. finer step size results in less degradation. 1% degradation seems adequate, which allows a representation of the self-apodised spectra in 16 bits per sample

# Possible Approach for IASI Data Compression

- Off-the-shelf loss-less compression is ineffective
  - Unix utility gzip provides 7.5% data reduction on a IASI Level 1c spectrum
- The data compression must be adapted to the characteristics of the IASI spectra; necessary steps include
  - Suitable representation of the spectra
  - Carefully controlled quantisation
  - Entropy encoding
  - Reverse procedure at user end

# Representation by EOF Decomposition

- Due to the atmospheric gas absorption all IASI spectra have a characteristic, similar shape
- EOF decomposition, based on pre-calculated eigenvectors, allows an efficient representation of the spectra by the eigenvector scores
- A lossy compression can be achieved by restricting the number of scores to the most pronounced co-variation structures, which is achieved by
  - Ranking the eigenvectors according to their eigenvalues (early-ranking, high-value, eigenvectors explain meteorological variability, low-ranking, low value, ones mainly represent noise)
  - Truncation of the representation of spectra by using only early-ranking, high eigenvalue, eigenvectors



# Encoding of Spectrum

- **Offline: definition of constants**
  - $NE\Delta R$ : Noise radiance spectrum
  - $U$ : Matrix of  $k$  (Order 100-300) column eigenvectors, describing ensemble of training set  $NE\Delta R$  normalised spectra
  - Huffman Code
  - Minor constants, step sizes, data boundaries etc.
- **Online: EOF scores**
  - Starting with individual spectrum  $y'$
  - Noise-normalise to  $y = y' / NE\Delta R$
  - Calculate scores  $c = U^T y$
  - Quantise according to pre-defined step-size to integer vector  $c'$
  - Use Huffman code to substitute  $c'$  integers by bit-stream

# Encoding of Spectrum (cont.)

- **Online: Residuals**
  - Subtract quantised-scores derived normalised spectrum from  $y$  to produce residual  $\Delta y = y - U c'$
  - Use Huffman code to represent quantised  $\Delta y$  as bit-stream
- **Online: Communication**
  - The bit-streams are communicated
- **Online: Decoding**
  - Receive bit-stream and decode to EOF scores  $c'$  and residuals  $\Delta y$
  - Use EOF scores  $c'$ ,
  - or reconstruct and use truncated spectrum  $\hat{y} = U c'$ ,
  - or add  $\Delta y$  to reconstruct and use un-truncated spectrum
  - Modified  $U$  can incorporate user-required spectral manipulation at no extra cost

# Truncated Spectra and Noise Reduction

- A decode scheme of great potential value for routine operations is to avoid full decode to observed spectra, but to synthesise the spectra from the truncated set of EOF scores only  $\hat{\mathbf{y}} = \mathbf{U}\mathbf{c}'$
- Smith and Woolf (1976):
  - Effects of random observation errors are minimised without suppressing real information
- Huang and Antonelli (2001):
  - Interferometric measurements based on 3888 samples can be represented, without loss within small fraction of measurement noise, by 150 scores
- Noise is reduced in RMS terms via an implicit filter based on past experience (eigenvector training set) and the implied decision to ignore the potential value of eigenvectors below truncation cut-off
- **Risk:** Filtering of observations based on past experience may well suppress reporting of unexpected conditions - the reason observations are made

# Caveats

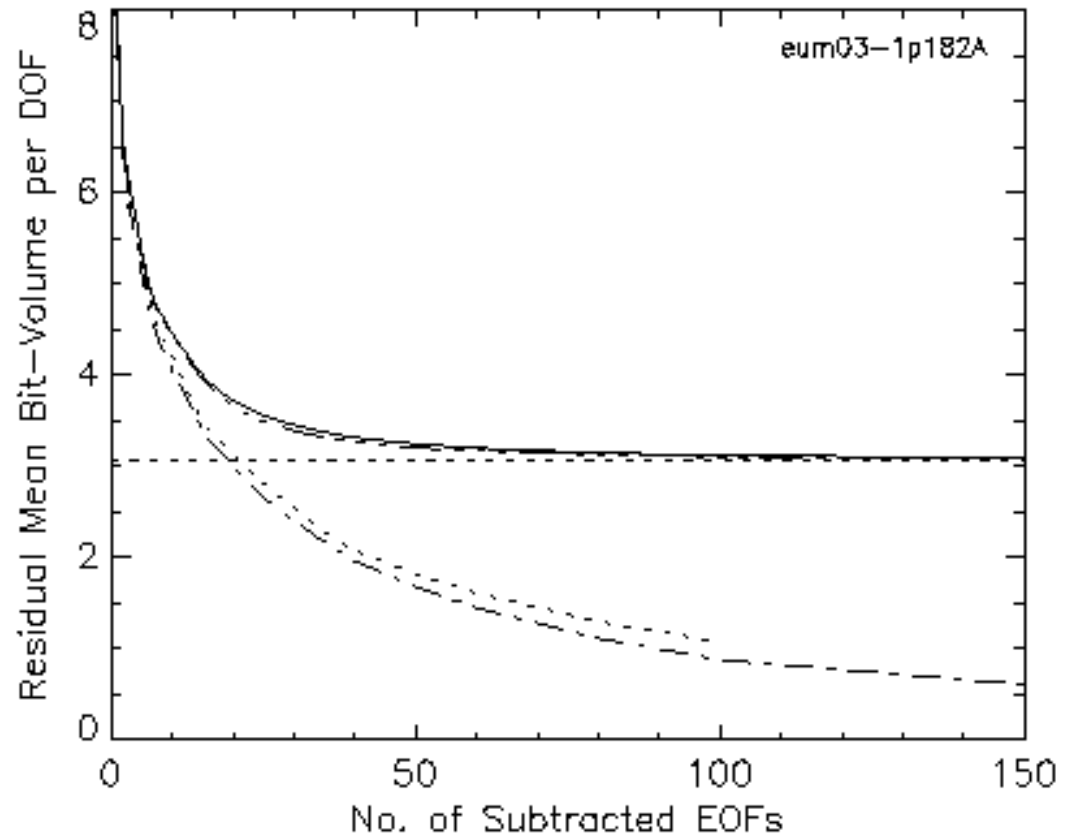
- **Correlation between adjacent channels**
  - The apodised spectra have non-diagonal error covariance
  - Noise of the apodised spectra is not white
- **Way forward**
  - Start with standardised self-apodised spectra, having diagonal error covariance and white noise characteristics
  - Provide user with post-processing software doing the apodisation, if desired

# Impact of Apodisation

- The concept of  $NE\Delta T$  is strictly only meaningful for self-apodised IASI spectra where “channels” are independent
- Users prefer heavily apodised spectra, where the self-apodised interferogram is attenuated by a factor  $\sim 30$  near OPD limits, compared to the one near zero OPD, resulting in gain attenuation of higher spatial frequencies
- The quantisation of an apodised spectrum bears the danger of swamping the attenuated higher frequencies by unattenuated higher frequencies of the quantisation noise, which can only be avoided by finer quantisation
- The full information contained in the IASI spectra will be retained only if the apodised spectra are quantised at 22 bit per sample

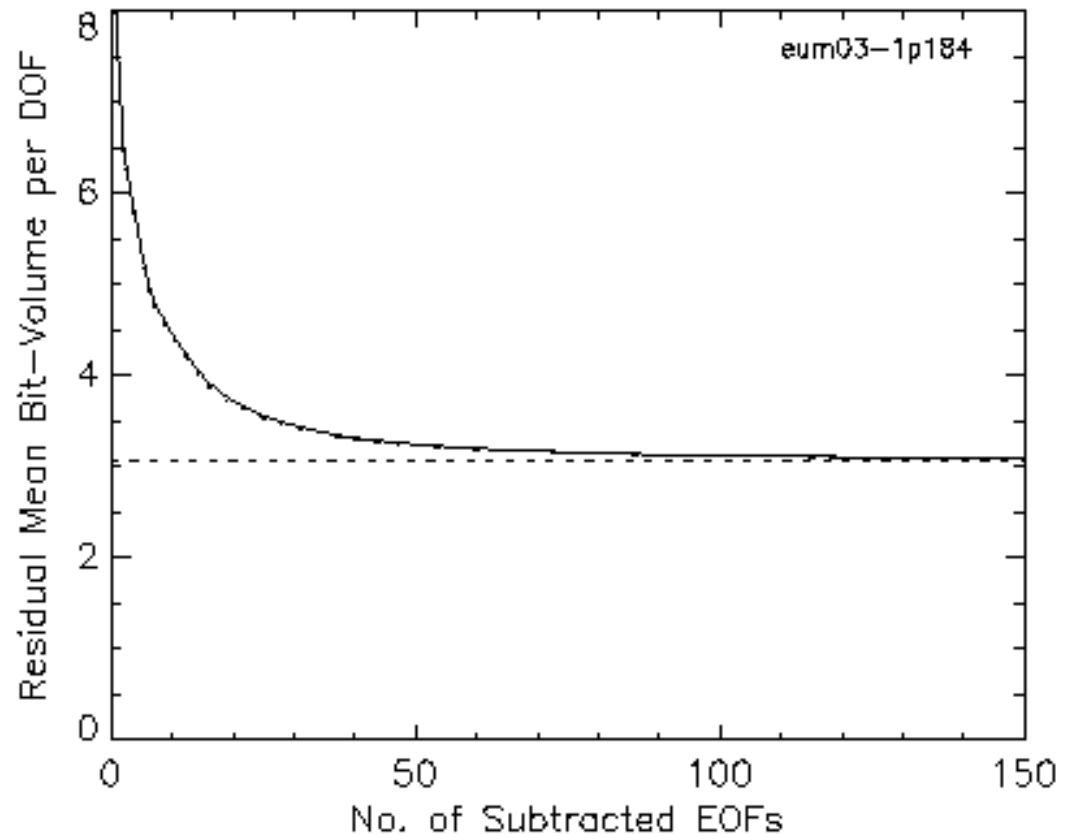
# Data Volume - PC Scores Only: Clear-Air Data

- Based on RTIASI-3.2 simulations using sub-sampled Chevallier data set: 1000 cases to train and 2373 cases to test
- Use noise-normalised radiance spectra:
  - In this case the “noise” is simulated
  - Actual noise may be included, or excluded
- For zero subtracted EOFs the “residual” entropy is large, about 9 bit  $\text{DOF}^{-1}$  must be accommodated
- For spectra containing noise the residual entropy falls to an asymptotic limit around 3 bit  $\text{DOF}^{-1}$  somewhere between 50 and 100 EOFs
- Where no actual noise included, the residual entropy falls to 3 bit  $\text{DOF}^{-1}$  at about 20 EOFs, and continues to fall



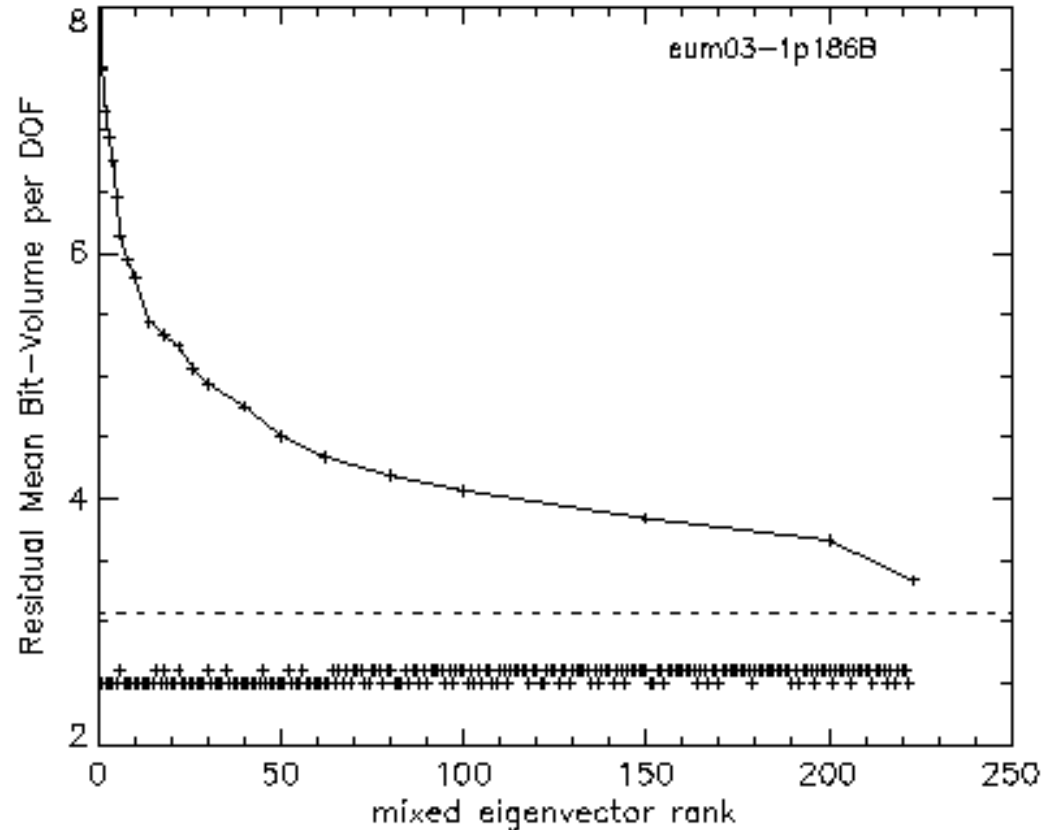
# Data Volume - PC Scores Only: Clear-Air Data Independent Test with Perturbed Profiles

- 1000 test profiles with added temperature perturbations
- Mean strength (single-Dirac) or dipole-moment (double-Dirac) of 1 km 1K
- Result: negligible difference in residual entropy at given number of EOFs
- 100 EOF scores are used to represent the clear-air IASI spectra



# Data Volume - PC Scores Only: Cloudy Profiles

- Based on von Bremen cloudy spectra (RTIASI-3.2 plus cloud parameterisation): 4000 sub-sampled, two thirds for training, one third for testing
- Clear-air properties are scraped out using the 100 clear-air EOFs, leaving a set of residual spectra not usually represented in clear air
- Normalised SVD is applied to give cloud-signature eigenvectors
- Cloud-signature EOFs are added to the clear-air EOFs such that the highest ranking eigenvalues are similar for both types
- The residual is about  $3.41 \text{ bit DOF}^{-1}$  for the cloudy spectra





# Residual Encoding

- The **quantised** EOF scores will be used to derive the residuals
- The residuals themselves are quantised at a step size of half- $NE\Delta R$ , thus increasing the  $NE\Delta R$  by no more than 1%
- The resulting integer amplitudes are Huffman encoded, based on their probability distribution, requiring 41 different descriptors
- The encoded residuals require 3.25 bits sample<sup>-1</sup> or 27500 bits spectrum<sup>-1</sup>

# Data Volumes

Inter-comparison of full and reduced data volumes

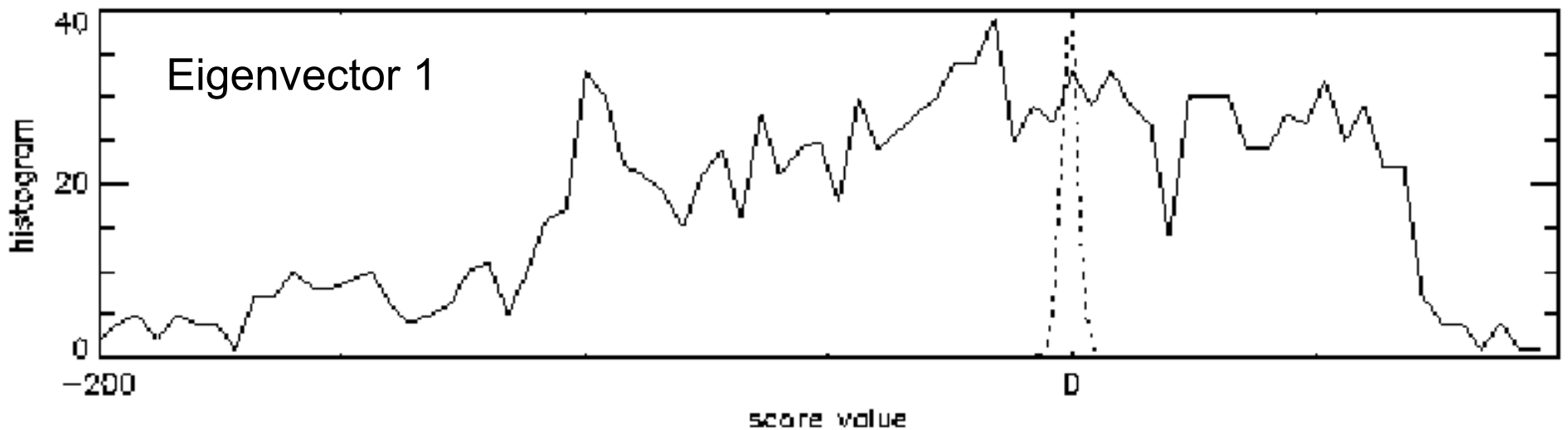
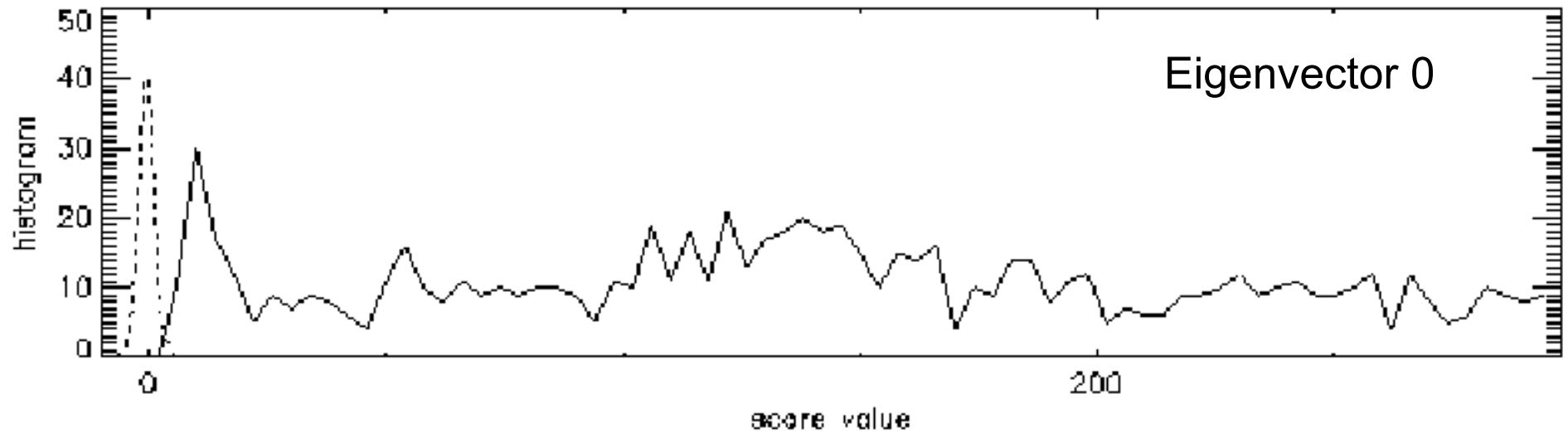
<b>IASI Data Representation</b>	<b>Data Volume (MB per scan line)</b>	<b>Data Volume (GB per day)</b>
Full spectrum, 24 bits sample <sup>-1</sup> (loss-free)	3.05	32.9
Full spectrum, 16 bits sample <sup>-1</sup> (slightly lossy)	2.03	21.9
300 selected channels, 24 bits sample <sup>-1</sup>	0.11	1.17
200 PC scores, 24 bits sample <sup>-1</sup> (lossy)	0.073	0.780
200 PC scores, compressed (lossy)	0.039	0.415
200 PC scores + residuals, compressed (loss-free)	0.45	2.59

# Data Processing

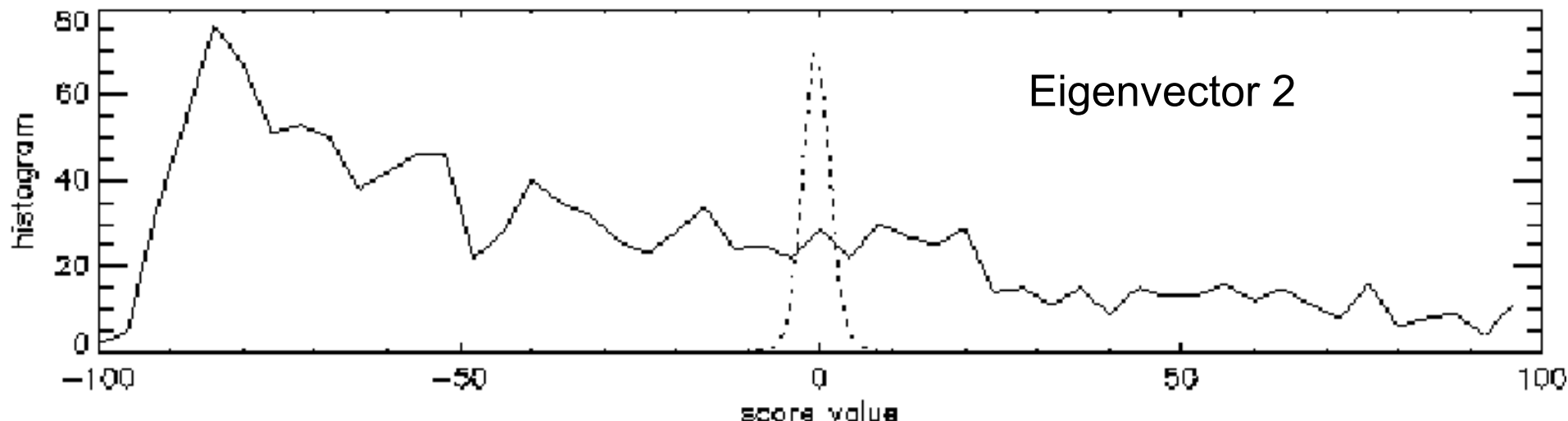
Transformation of IASI spectra into PC scores and compression requires additional processing steps at CGS and user ends

- Offline definition of constants and code
- **CGS processing:** Noise-normalise self-apodised spectra, calculation of PC scores and residuals, quantisation into integers, run-length and Huffman encoding, communication of bit-stream
- **Processing at user end:** Receipt of bit-stream, reverse processing to decode into PC scores and residuals, reconstruction of spectra from PC scores, optional adding of residuals, apodisation of spectra; likewise direct use of PC scores for cloud detection, quality check, geophysical parameters retrieval, or assimilation
- Despite all this, processing and retrieval stages can involve less CPU resources than needed for processing of uncoded spectra

# Cloud Detection Using EOF Scores



## Cloud Detection Using EOF Scores (cont.)



- Cloud detection using thresholds on scores of first few eigenvectors allows for efficient cloud detection
- Using 10 cloud-signature eigenvectors produces no false decisions over all clear and cloudy spectra
- Choice of threshold is uncritical, any value between 25 and 65 is fine
- But: This “perfect” result probably because “von Bremen Cloud” under-represents marginal cloud situations

# Representation of FRTM in EOF Space

- Fast Radiative Transfer Models (FRTM), including their adjoint and tangent-linear versions, are essential tools to assimilate satellite-measured radiances
- FRTMs have been developed for the use with hyper-spectral sounder data as well (e.g. RTIASI, RTTOV, SARTA), but the huge number of spectral samples/channels prevent the assimilation of full spectra
- In view of the possibility to represent the (almost) full information of the hyperspectral sounders in few hundred EOF scores it seems appropriate to seek for a representation of the FRTM in EOF space
- EOF scores are linear combinations of radiance samples
  - The scores can be considered as convolution of monochromatic radiances with a “response function” that is described by the eigenvectors

# FRTM in EOF Space

Courtesy: X. Liu (unpublished material)

Radiance spectrum:

$$\mathbf{y} = \mathbf{U}\mathbf{c} = \sum_{i=1}^{N_{EOF}} c_i \mathbf{u}_i + \boldsymbol{\varepsilon}$$

EOF scores:

$$c_i = \sum_{j=1}^{N_{Ch}} u_{ij} y_j$$

Convolved radiance:

$$y_j = \sum_{k=1}^{N_{mono}} a_{jk} y_k^{mono}$$

# FRTM in EOF Space (cont.)

Courtesy: X. Liu (unpublished material)

Convolved radiance:

$$y_j = \sum_{k=1}^{N_{mono}} a_{jk} y_k^{mono}$$

$\Rightarrow$

$$c_i = \sum_{j=1}^{N_{Ch}} u_{ji} \sum_{k=1}^{N_{mono}} a_{kj} y_k$$

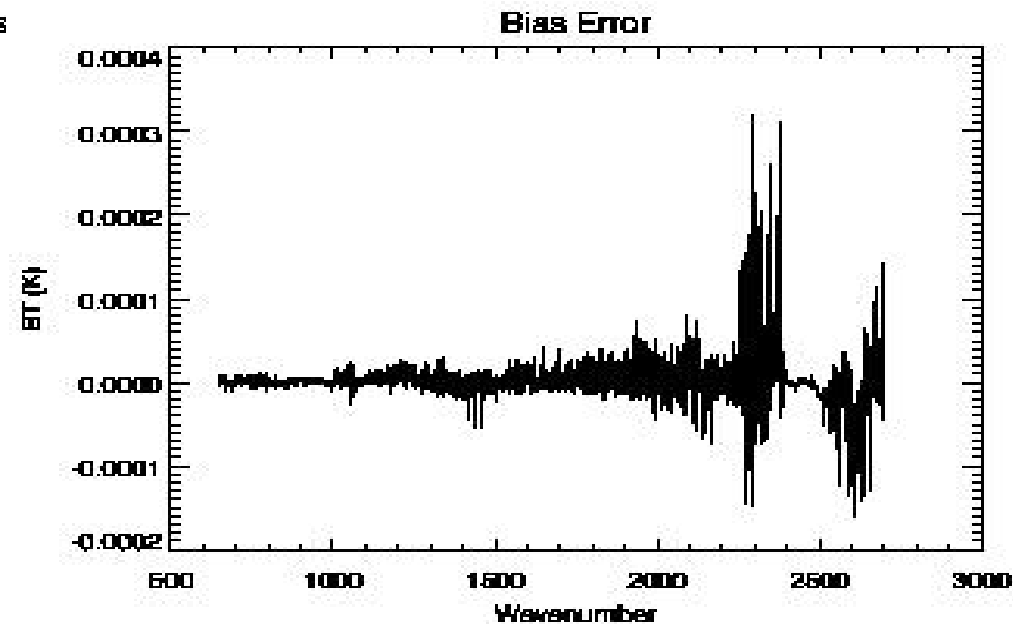
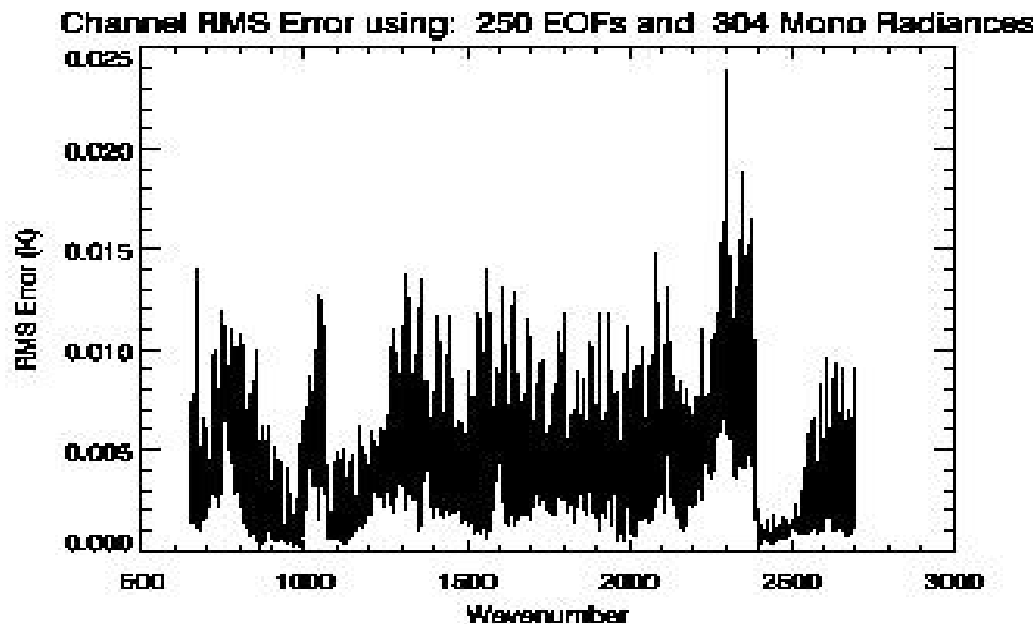
$$c_i = \sum_{k=1}^{N_{mono}} y_k^{mono} \sum_{j=1}^{N_{Ch}} u_{ji} a_{kj} = \sum_{k=1}^{N_{mono}} y_k^{mono} b_{ki}$$



# FRTM in EOF Space (cont.)

## Courtesy: X. Liu (unpublished material)

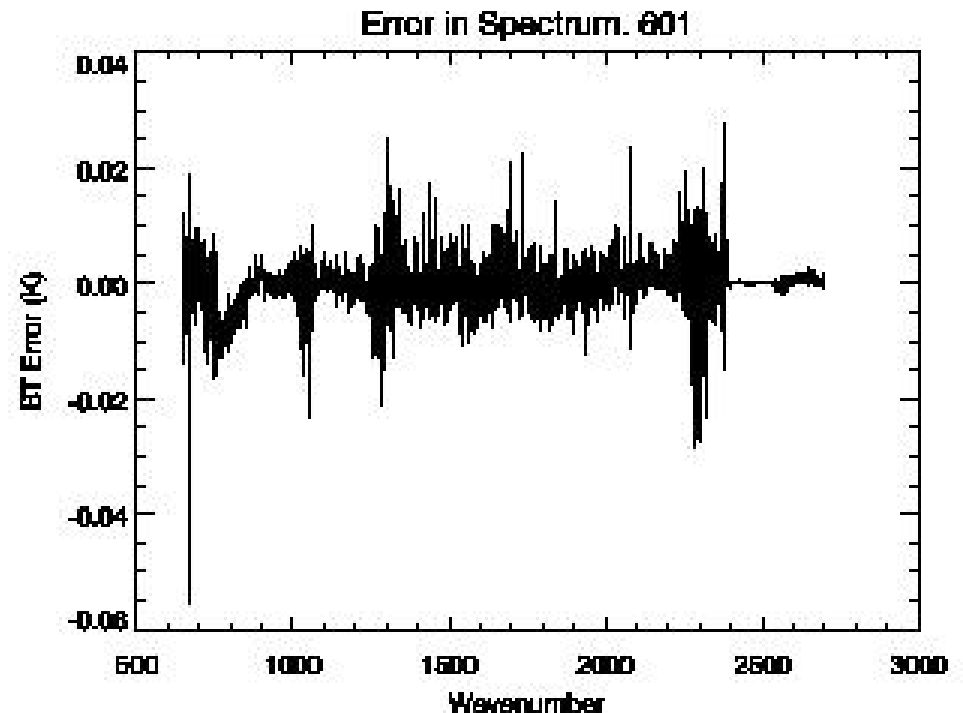
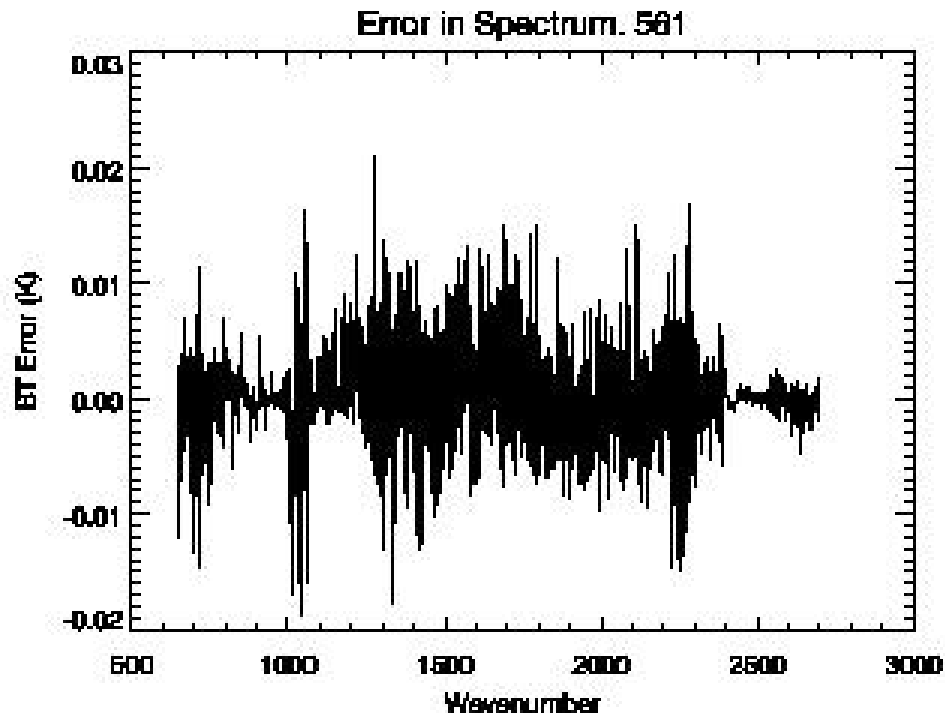
- Training of the model for the NAST-I interferometer (8632 spectral samples)
- RMS errors within 0.025 K brightness temperature
- Biases within (-0.0002 K, 0.0004 K) brightness temperature



# FRTM in EOF Space (cont.)

Courtesy: X. Liu (unpublished material)

Typical errors in selected spectra rarely exceed 0.05 K

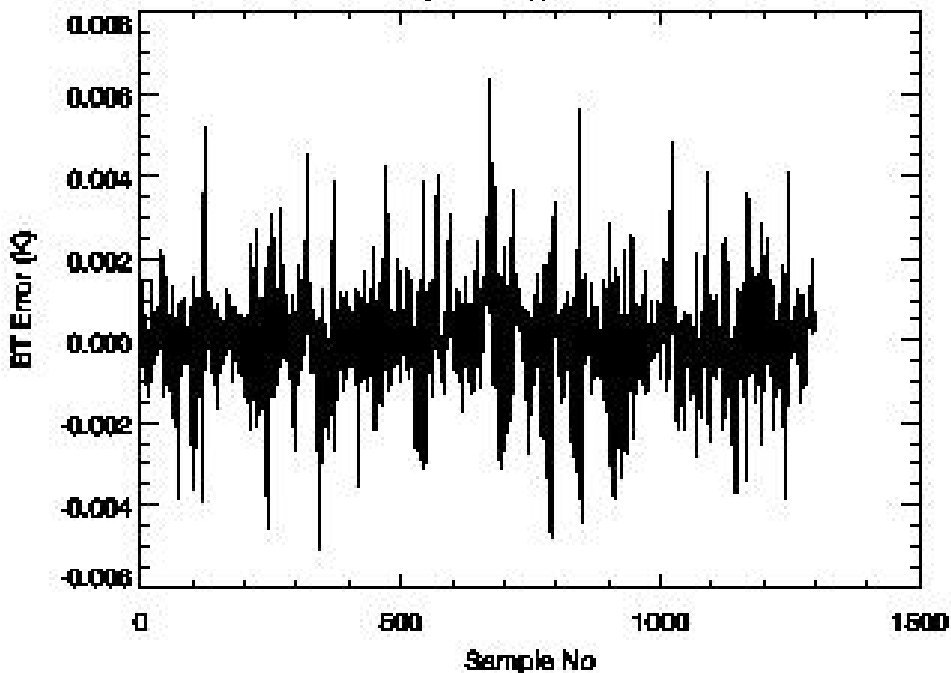


# FRTM in EOF Space (cont.)

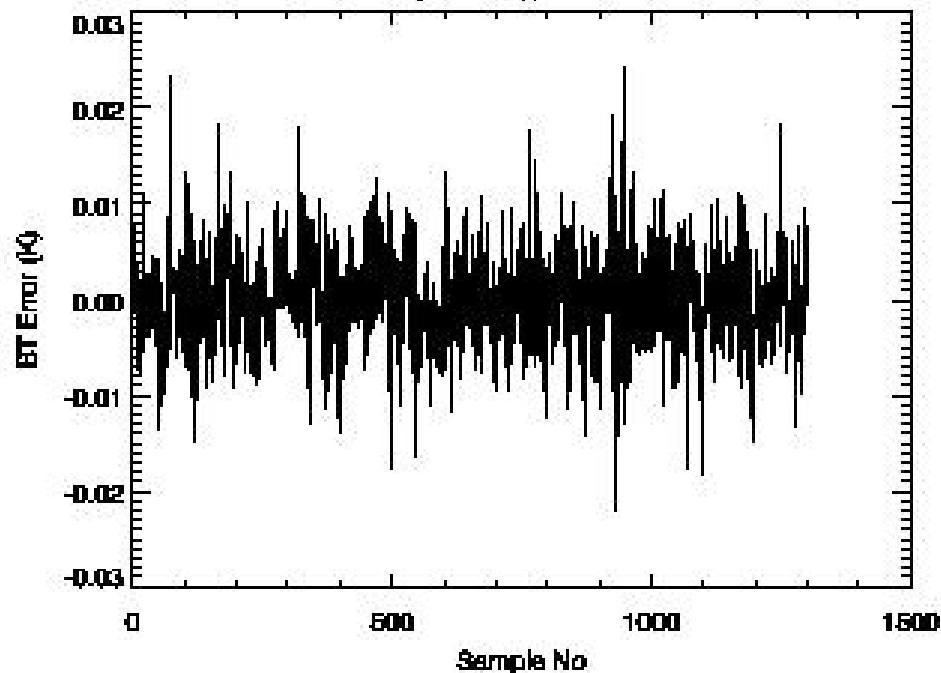
## Courtesy: X. Liu (unpublished material)

Errors for selected spectral samples at  $1127.3 \text{ cm}^{-1}$  and  $1309.6 \text{ cm}^{-1}$  covering 1300 different situations

Ch2001 (1127.3), RMS = 0.0014



Ch2801 (1309.6), RMS = 0.0055



# FRTM in EOF Space (cont.)

## Courtesy: X Liu (unpublished material)

- Detailed numerical study shows that the number of monochromatic radiances needed to provide a reasonable fit is of the order 300
- The spectra are re-constructed using 250 EOF scores
- The RMS errors are less than 0.025 K in brightness temperature for all spectral samples
- Considering forward calculations for the entire spectrum the proposed model is one to two orders of magnitude faster than other fast models, so that it offers the possibility to assimilate the entire spectral information

# Conclusion

- IASI spectra can be efficiently compressed by EOF decomposition followed by quantisation and subsequent entropy encoding
- Loss-less data-volume reduction at a factor  $\sim 7$  is achievable
- Acceptance of moderate information loss can **further** reduce the volume substantially, by an additional factor of 20-80
- The representation of the IASI spectra in terms of eigenvector scores is of direct benefit for efficient cloud detection
- Super-fast RT models are being developed that simulate eigenvector scores instead of radiance spectra
- Assimilation of leading eigenvector scores will add more information to NWP than single “channels” and fully exploit the information contained in the IASI spectra