

Handling Provenance, Including Mathematical Proofs, in Cyberinfrastructure- Oriented Data Processing

Ann Q. Gates¹, Olga Kosheleva¹,
Vladik Kreinovich¹, Sa-aat Niwitpong²,
Paulo Pinheiro da Silva¹, and Leonardo Salayandia¹

¹University of Texas, El Paso, Texas 79968, USA
contact email vladik@utep.edu

²King Mongkut's University of Technology North Bangkok

Cyberinfrastructure: . . .

Data Processing vs. . . .

Need for Uncertainty . . .

Case Study: Seismic . . .

Proofs as a Particular . . .

Conclusions

Acknowledgments

Title Page



Page 1 of 20

Go Back

Full Screen

Close

Quit

1. Cyberinfrastructure: A Brief Overview

- Practical problem: need to combine geographically separate computational resources.
- Centralization of computational resources – traditional approach to combining computational resources.
- Limitations of centralization:
 - need to reformat all the data;
 - need to rewrite data processing programs: make compatible w/selected formats and w/each other
- Cyberinfrastructure – a more efficient approach to combining computational resources:
 - keep resources at their current locations, and
 - in their current formats.
- Technical advantages of cyberinfrastructure: a brief summary.

Cyberinfrastructure: ...

Data Processing vs. ...

Need for Uncertainty ...

Case Study: Seismic ...

Proofs as a Particular ...

Conclusions

Acknowledgments

Title Page



Page 2 of 20

Go Back

Full Screen

Close

Quit

2. Data Processing vs. Data Fusion

- *Practically important situation:* difficult to measure the desired quantity y with a given accuracy.
- *Data processing:*
 - measure related easier-to-measure quantities x_1, \dots, x_n ;
 - estimate y from the results \tilde{x}_i of measuring x_i as $\tilde{y} = f(\tilde{x}_1, \dots, \tilde{x}_n)$.
- *Example:* seismic inverse problem.
- *Data fusion:*
 - measure the quantity y several times;
 - combine the results $\tilde{y}_1, \dots, \tilde{y}_n$ of these measurements.
- *Specifics of cyberinfrastructure:* first looks for stored results \tilde{x}_i (corr., \tilde{y}_i), measure only if necessary.
- *Combination of data processing and data fusion.*

Cyberinfrastructure: ...

Data Processing vs. ...

Need for Uncertainty ...

Case Study: Seismic ...

Proofs as a Particular ...

Conclusions

Acknowledgments

Title Page



Page 3 of 20

Go Back

Full Screen

Close

Quit

3. Need for Uncertainty Propagation, and for Provenance of Uncertainty

- *Need for uncertainty propagation.*
 - main reasons for data processing and data fusion: accuracy is not high enough;
 - we must make sure that after the data processing (data fusion), we get the desired accuracy.
- *In cyberinfrastructure this is especially important:*
 - accuracy varies greatly, and
 - we do not have much control over these accuracies.
- *Need for the provenance of uncertainty:*
 - sometimes, the resulting accuracy is still too low;
 - it is desirable to find out which data points contributed most to the inaccuracy.

Cyberinfrastructure: ...

Data Processing vs. ...

Need for Uncertainty ...

Case Study: Seismic ...

Proofs as a Particular ...

Conclusions

Acknowledgments

Title Page



Page 4 of 20

Go Back

Full Screen

Close

Quit

4. Uncertainty of the Results of Direct Measurements: Probabilistic and Interval Approaches

- Manufacturer of the measuring instrument (MI) supplies Δ_i s.t. $|\Delta x_i| \leq \Delta_i$, where $\Delta x_i \stackrel{\text{def}}{=} \tilde{x}_i - x_i$.
- The actual (unknown) value x_i of the measured quantity is in the interval $\mathbf{x}_i = [\tilde{x}_i - \Delta_i, \tilde{x}_i + \Delta_i]$.
- *Probabilistic uncertainty*: often, we know the probabilities of different values $\Delta x_i \in [-\Delta_i, \Delta_i]$.
- *How probabilities are determined*: by comparing our MI with a much more accurate (standard) MI.
- *Interval uncertainty*: in two cases, we do not determine the probabilities:
 - cutting-edge measurements;
 - measurements on the shop floor.
- In both cases, we only know that $x_i \in [\tilde{x}_i - \Delta_i, \tilde{x}_i + \Delta_i]$.

Cyberinfrastructure: ...

Data Processing vs. ...

Need for Uncertainty ...

Case Study: Seismic ...

Proofs as a Particular ...

Conclusions

Acknowledgments

Title Page

◀◀ ▶▶

◀ ▶

Page 5 of 20

Go Back

Full Screen

Close

Quit

5. Typical Situation: Measurement Errors are Reasonably Small

- *Typical situation:*
 - direct measurements are accurate enough;
 - the resulting approximation errors Δx_i are small;
 - terms which are quadratic (or of higher order) in Δx_i can be safely neglected.
- *Example:* for an error of 1%, its square is a negligible 0.01%.
- *Linearization:*
 - expand f in Taylor series around the point $(\tilde{x}_1, \dots, \tilde{x}_n)$;
 - restrict ourselves only to linear terms:

$$\Delta y = c_1 \cdot \Delta x_1 + \dots + c_n \cdot \Delta x_n,$$

where $c_i \stackrel{\text{def}}{=} \frac{\partial f}{\partial x_i}$.

Cyberinfrastructure: . . .

Data Processing vs. . . .

Need for Uncertainty . . .

Case Study: Seismic . . .

Proofs as a Particular . . .

Conclusions

Acknowledgments

Title Page



Page 6 of 20

Go Back

Full Screen

Close

Quit

6. Case of Data Processing

- *Propagation (probabilistic case)*: if Δx_i are independent with st. dev. σ_i (and 0 mean), then Δy has st. dev.

$$\sigma^2 = c_1^2 \cdot \sigma_1^2 + \dots + c_n^2 \cdot \sigma_n^2.$$

- *Provenance*:
 - we know which component σ^2 comes from the i -th measurement;
 - we can predict how replacing the i -th measurement with a more accurate one ($\sigma_i^{\text{new}} \ll \sigma_i$) will affect σ^2 .
- *Propagation of interval uncertainty*:

$$\Delta = |c_1| \cdot \Delta_1 + \dots + |c_n| \cdot \Delta_n.$$

- We can predict how replacing the i -th measurement with a more accurate one ($\Delta_i^{\text{new}} \ll \Delta_i$) will affect Δ .

Cyberinfrastructure: ...

Data Processing vs. ...

Need for Uncertainty ...

Case Study: Seismic ...

Proofs as a Particular ...

Conclusions

Acknowledgments

Title Page



Page 7 of 20

Go Back

Full Screen

Close

Quit

7. Beyond Probabilistic and Interval Uncertainty

- *Up to now*: we considered two extreme situations:
 - *probabilistic* uncertainty, when we know all the probabilities;
 - *interval* uncertainty, when we have no information about the probabilities.
- *Fact*: probabilistic situation is a particular case of the interval situation.
- *Conclusion*: interval bounds are wider.
- *In practice*: often, we have partial information about probabilities.
- *As a result*:
 - probabilistic bounds are too narrow,
 - interval bounds are too wide.
- *We need*: intermediate bounds.

Cyberinfrastructure: . . .

Data Processing vs. . . .

Need for Uncertainty . . .

Case Study: Seismic . . .

Proofs as a Particular . . .

Conclusions

Acknowledgments

Title Page



Page 8 of 20

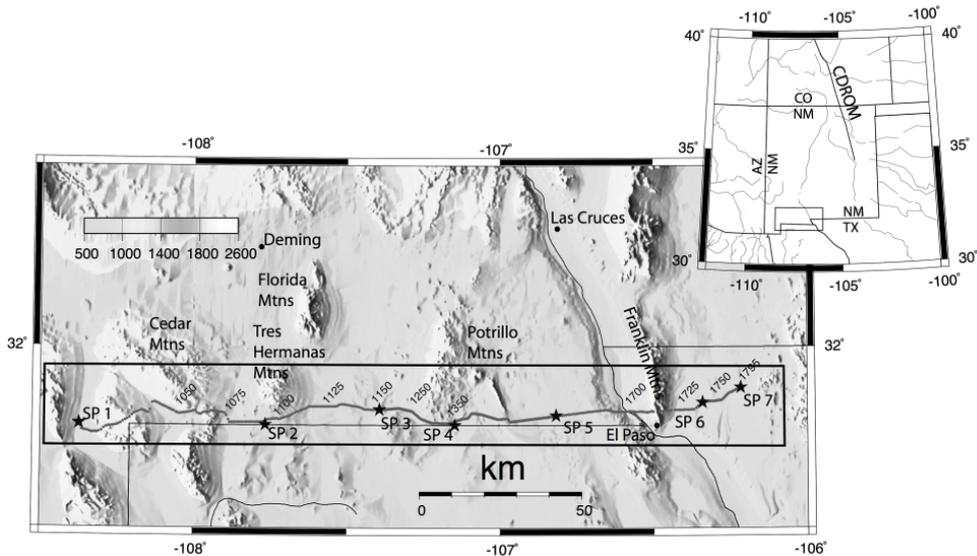
Go Back

Full Screen

Close

Quit

8. Case Study: Seismic Inverse Problem in the Geosciences



Cyberinfrastructure: ...

Data Processing vs. ...

Need for Uncertainty ...

Case Study: Seismic ...

Proofs as a Particular ...

Conclusions

Acknowledgments

Title Page

◀ ▶

◀ ▶

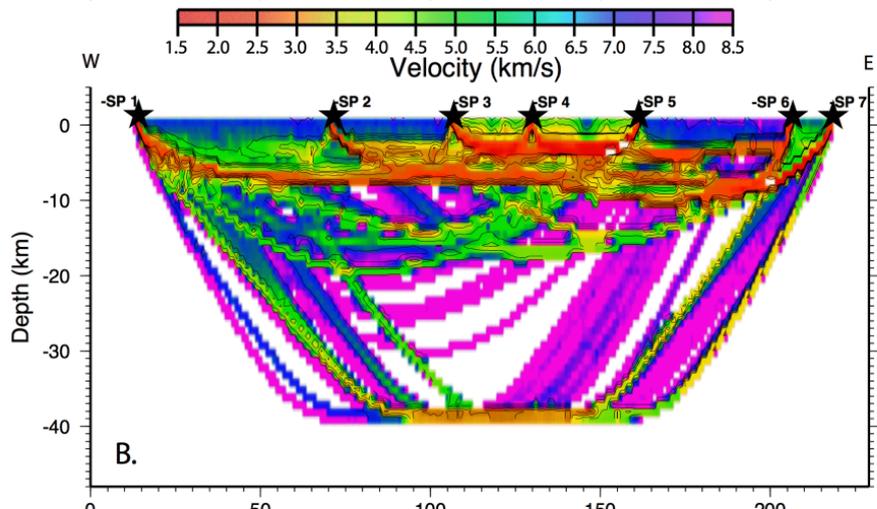
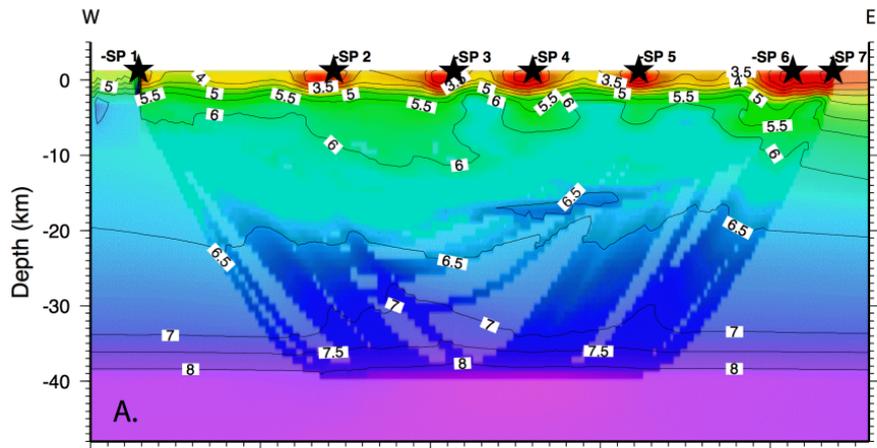
Page 9 of 20

Go Back

Full Screen

Close

Quit



Cyberinfrastructure: ...

Data Processing vs. ...

Need for Uncertainty ...

Case Study: Seismic ...

Proofs as a Particular ...

Conclusions

Acknowledgments

Title Page

◀ ▶

◀ ▶

Page 10 of 20

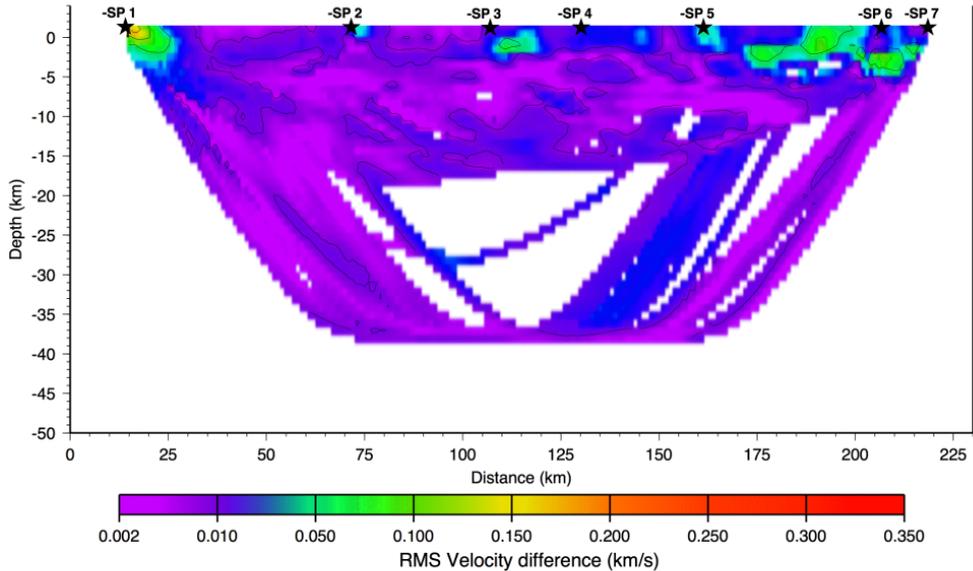
Go Back

Full Screen

Close

Quit

9. Estimating Uncertainty, First Try: Probabilistic Approach



Cyberinfrastructure: ...

Data Processing vs. ...

Need for Uncertainty ...

Case Study: Seismic ...

Proofs as a Particular ...

Conclusions

Acknowledgments

Title Page

◀ ▶

◀ ▶

Page 11 of 20

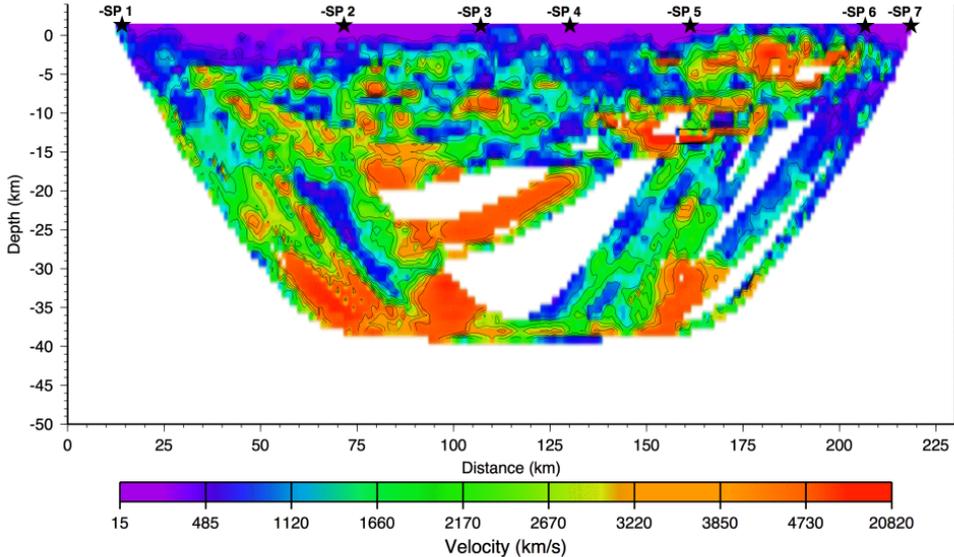
Go Back

Full Screen

Close

Quit

10. Estimating Uncertainty, Second Try: Interval Approach



Cyberinfrastructure: ...

Data Processing vs. ...

Need for Uncertainty ...

Case Study: Seismic ...

Proofs as a Particular ...

Conclusions

Acknowledgments

Title Page

◀◀ ▶▶

◀ ▶

Page 12 of 20

Go Back

Full Screen

Close

Quit

11. Towards a Better Estimate: Revisiting Estimation Algorithms Under Probabilistic and Interval Uncertainty

- *Linearization*: $\Delta y = \sum_{i=1}^n c_i \cdot \Delta x_i$, where $c_i \stackrel{\text{def}}{=} \frac{\partial f}{\partial x_i}$.
- *Formulas*: $\sigma^2 = \sum_{i=1}^n c_i^2 \cdot \sigma_i^2$, $\Delta = \sum_{i=1}^n |c_i| \cdot \Delta_i$.
- *Numerical differentiation*: n iterations, too long.
- *Monte-Carlo approach*: if Δx_i are Gaussian w/ σ_i , then $\Delta y = \sum_{i=1}^n c_i \cdot \Delta x_i$ is also Gaussian, w/desired σ .
- *Advantage*: # of iterations does not grow with n .
- *Interval estimates*: if Δx_i are Cauchy, w/ $\rho_i(x) = \frac{\Delta_i}{\Delta_i^2 + x^2}$, then $\Delta y = \sum_{i=1}^n c_i \cdot \Delta x_i$ is also Cauchy, w/desired Δ .

Cyberinfrastructure: ...

Data Processing vs. ...

Need for Uncertainty ...

Case Study: Seismic ...

Proofs as a Particular ...

Conclusions

Acknowledgments

Title Page



Page 13 of 20

Go Back

Full Screen

Close

Quit

12. Resulting Fast (Linearized) Algorithm for Estimating Interval Uncertainty

- Apply f to \tilde{x}_i : $\tilde{y} := f(\tilde{x}_1, \dots, \tilde{x}_n)$;
- For $k = 1, 2, \dots, N$, repeat the following:
 - use RNG to get $r_i^{(k)}$, $i = 1, \dots, n$ from $U[0, 1]$;
 - get st. Cauchy values $c_i^{(k)} := \tan(\pi \cdot (r_i^{(k)} - 0.5))$;
 - compute $K := \max_i |c_i^{(k)}|$ (to stay in linearized area);
 - simulate “actual values” $x_i^{(k)} := \tilde{x}_i - \delta_i^{(k)}$, where $\delta_i^{(k)} := \Delta_i \cdot c_i^{(k)} / K$;
 - simulate error of the indirect measurement:
$$\delta^{(k)} := K \cdot \left(\tilde{y} - f \left(x_1^{(k)}, \dots, x_n^{(k)} \right) \right);$$
- Solve the ML equation $\sum_{k=1}^N \frac{1}{1 + \left(\frac{\delta^{(k)}}{\Delta} \right)^2} = \frac{N}{2}$ by bisection, and get the desired Δ .

Cyberinfrastructure: ...

Data Processing vs. ...

Need for Uncertainty ...

Case Study: Seismic ...

Proofs as a Particular ...

Conclusions

Acknowledgments

Title Page

◀ ▶

◀ ▶

Page 14 of 20

Go Back

Full Screen

Close

Quit

13. A New (Heuristic) Approach

- *Problem:* guaranteed (interval) bounds are too high.
- *Gaussian case:* we only have bounds guaranteed with confidence, say, 90%.
- *How:* cut top 5% and low 5% off a normal distribution.
- *New idea:* to get similarly estimates for intervals, we “cut off” top 5% and low 5% of Cauchy distribution.
- *How:*
 - find the threshold value x_0 for which the probability of exceeding this value is, say, 5%;
 - replace values x for which $x > x_0$ with x_0 ;
 - replace values x for which $x < -x_0$ with $-x_0$;
 - use this “cut-off” Cauchy in error estimation.
- *Example:* for 95% confidence level, we need $x_0 = 12.706$.

Cyberinfrastructure: ...

Data Processing vs. ...

Need for Uncertainty ...

Case Study: Seismic ...

Proofs as a Particular ...

Conclusions

Acknowledgments

Title Page



Page 15 of 20

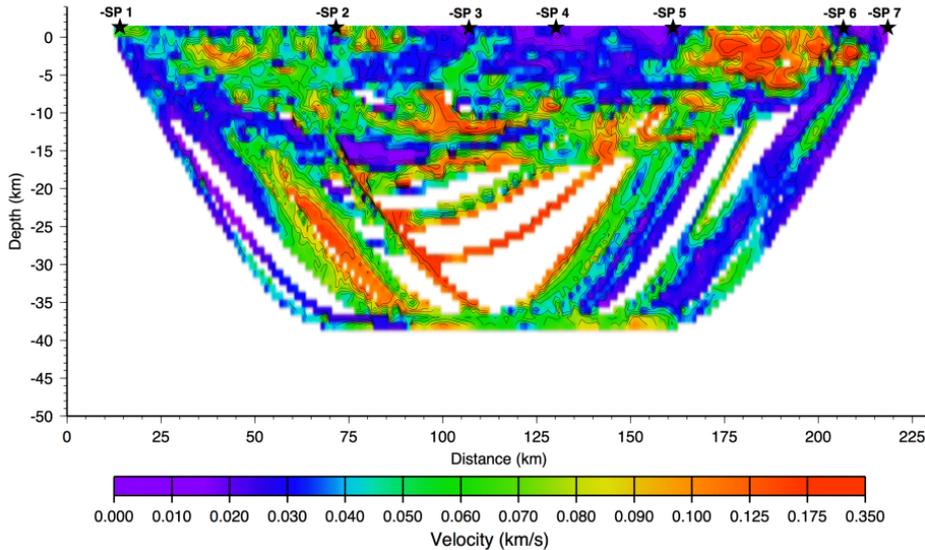
Go Back

Full Screen

Close

Quit

14. Heuristic Approach: Results with 95% Confidence Level



Cyberinfrastructure: ...

Data Processing vs. ...

Need for Uncertainty ...

Case Study: Seismic ...

Proofs as a Particular ...

Conclusions

Acknowledgments

Title Page



Page 16 of 20

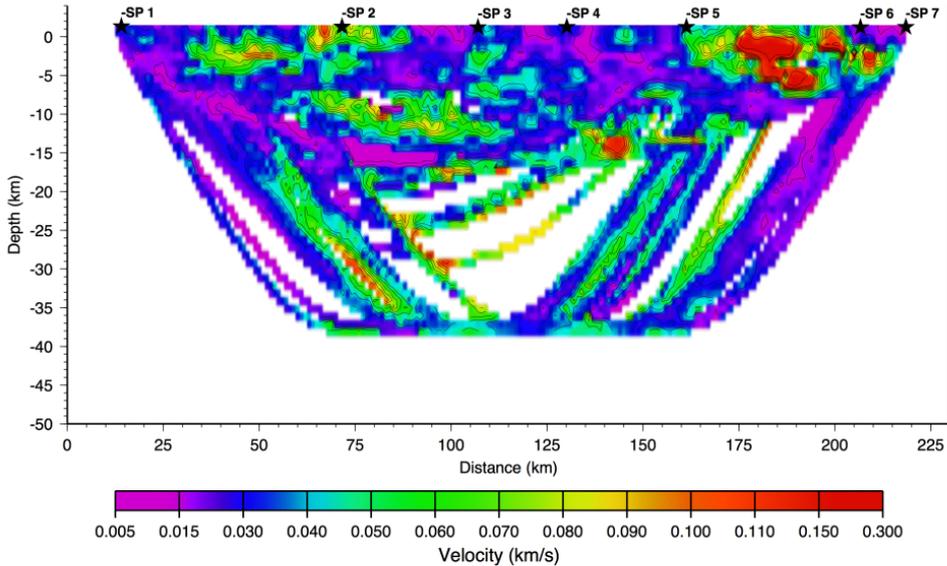
Go Back

Full Screen

Close

Quit

15. Heuristic Approach: Results with 90% Confidence Level



Cyberinfrastructure: ...

Data Processing vs. ...

Need for Uncertainty ...

Case Study: Seismic ...

Proofs as a Particular ...

Conclusions

Acknowledgments

Title Page



Page 17 of 20

Go Back

Full Screen

Close

Quit

16. Proofs as a Particular Case of Provenance

- For handling algorithms provenance, we need to handle different types of such provenance ranging:
 - from expert opinion on heuristic techniques
 - to experimental confirmation of semi-heuristic numerical methods
 - to formal proofs of algorithm correctness.
- In line with the main ideas behind cyberinfrastructure, it is desirable:
 - to combine and process these provenances
 - without moving them to a central location.
- This necessitates, e.g., a need to keep the proof of correctness of the combined algorithm de-centralized.

Cyberinfrastructure: ...

Data Processing vs. ...

Need for Uncertainty ...

Case Study: Seismic ...

Proofs as a Particular ...

Conclusions

Acknowledgments

Title Page



Page 18 of 20

Go Back

Full Screen

Close

Quit

17. Conclusions

- *In the past:* communications were much slower.
- *Conclusion:* use centralization.
- *At present:* communications are much faster.
- *Conclusion:* use cyberinfrastructure.
- *Related problems:*
 - gauge the the uncertainty of the results obtained by using cyberinfrastructure;
 - which data points contributed most to uncertainty;
 - how an improved accuracy of these data points will improve the accuracy of the result.
- *We described:* algorithms for solving these problems.
- *Additional problem:* what if interval estimates are too wide and probabilistic estimates are too narrow.

Cyberinfrastructure: . . .

Data Processing vs. . . .

Need for Uncertainty . . .

Case Study: Seismic . . .

Proofs as a Particular . . .

Conclusions

Acknowledgments

Title Page



Page 19 of 20

Go Back

Full Screen

Close

Quit

18. Acknowledgments

This work was supported in part:

- by NSF grant HRD-0734825 and
- by Grant 1 T36 GM078000-01 from the National Institutes of Health.

Cyberinfrastructure: . . .

Data Processing vs. . . .

Need for Uncertainty . . .

Case Study: Seismic . . .

Proofs as a Particular . . .

Conclusions

Acknowledgments

Title Page



Page 20 of 20

Go Back

Full Screen

Close

Quit