

Towards a Quantitative Cartography of the Grain Boundary Energy Landscape: Paths and Correlations: Supplementary Information

Sterling G. Baird^{a,*}, Eric R. Homer^a, David T. Fullwood^a, Oliver K. Johnson^a

^a*Department of Mechanical Engineering, Brigham Young University, Provo, UT 84602, USA*

Contents

| | |
|--|-----------|
| S1 Brief Summary of VFZ Methods | 1 |
| S2 Semivariograms for Estimating Global and Local Correlation Lengths | 2 |
| S2.1 Semivariogram Method | 2 |
| S2.2 Global Correlation Lengths | 3 |
| S2.3 Local Correlation Lengths | 4 |
| S3 Fe Input Dataset Characteristics | 4 |
| S3.1 Methods and Results | 7 |
| S3.2 Intrinsic Uncertainty | 7 |
| S3.3 Overprediction Bias | 8 |
| S3.4 Improving on Existing Datasets | 8 |
| S4 Gridded Sampling for Numerical Differentiation | 10 |
| S5 GBs Used for Path Visualization | 10 |
| Glossary | 10 |

S1. Brief Summary of VFZ Methods

We summarize the following aspects of the Voronoi fundamental zone (VFZ) framework:

- Creation and definition of a VFZ
- Mapping grain boundaries (GBs) to the VFZ
- Distance calculations
- Interpolation
- Comparison with traditional grain boundary octonion (GBO) metric

*Corresponding author.

Email address: ster.g.baird@gmail.com (Sterling G. Baird)

Each of these is described in greater detail in Baird et al. [1].

To define a VFZ, an arbitrary, fixed, low-symmetry reference GBO is chosen (o_{ref}) and for our use of GBOs, the VFZ is defined as the region of \mathbb{S}^7 (the unit 7-sphere in 8 dimensions) closer to o_{ref} than any of its symmetric images. If a low-symmetry GB is chosen, the point within a VFZ will be unique within numerical tolerance (and hence it is a true fundamental zone). Additionally, we use a Euclidean approximation to the true geodesic distance.

A GBO is composed of two quaternions, with the boundary plane normal in the $+z$ direction [2]. A GBO is mapped into a VFZ by calculating the pairwise distances between the reference GBO and each of the symmetrically equivalent octonions¹ and taking the symmetrically equivalent octonion closest to the reference GBO.

Once a GBO has been mapped into a VFZ, distance calculations proceed without further consideration of symmetrically equivalent octonions. The VFZ framework suffers from occasional, large distance overestimation which imposes a local sparseness of data and lead to poorer interpolation near the borders of a VFZ. However, this can be mitigated through ensemble or data augmentation techniques.

Note that in the original definition of the VFZ, the authors examined different interpolation techniques [1]. In the present work, we focus on Gaussian process regression which in our case imposes the assumption that crystallographically similar GBs share similar grain boundary energies (GBEs) within some correlation length. Gaussian process regression has the added benefit of built-in uncertainty quantification.

The primary differences between the VFZ framework and traditional GBO distance metric are that the VFZ framework is defined by a continuous set of points, exhibits occasional distance overestimation, uses a Euclidean approximation, and has a lower computational complexity.

S2. Semivariograms for Estimating Global and Local Correlation Lengths

S2.1. Semivariogram Method

First, we provide some background to put correlation lengths in context. The construction of Gaussian process regression models involves the combination of a prior distribution over the model space, with some set of observations and their quantified uncertainty. The result is a posterior distribution that provides the probability (density) of any particular model in light of the observed data and priors. In Gaussian process regression, as the name suggests, the priors are assumed to be Gaussian and therefore of the form

$$f(\mathbf{m}) \propto \exp \left(-\frac{1}{2}(\mathbf{m} - \mathbf{m}_0)^\top \mathbf{C}_{\mathbf{m}_0}^{-1}(\mathbf{m} - \mathbf{m}_0) \right) \quad (\text{S1})$$

where $f(\mathbf{m})$ is the probability density of an arbitrary model, $\mathbf{m} = \mathbf{m}(x)$, and $\mathbf{m}_0 = \mathbf{m}_0(x)$ is the prior model (i.e. a guess as to what the true model ought to look like). The quantity $\mathbf{C}_{\mathbf{m}_0} = \mathbf{C}_{\mathbf{m}_0}(x_i, x_j)$ is the “kernel” or covariance function of the prior and it describes the prior (assumed) covariance between the values $\mathbf{m}(x_i)$ and $\mathbf{m}(x_j)$.

It is possible to use a wide variety of kernel types, depending on the prior information one may have about the physical phenomenon one is attempting to model, e.g. continuity, differentiability, anisotropy, stationarity, and length-scales of correlation. One of the most common kernels employed is the Gaussian

¹Contrary to Francis et al. [2] which uses the passive convention for misorientation, we employ the active convention [1].

(sometimes called the squared exponential) kernel:

$$\mathbf{C}_{\mathbf{m}_0}(x_i, x_j) = \sigma^2 e^{-\frac{(x_i - x_j)^\top (x_i - x_j)}{2l^2}} \quad (\text{S2})$$

where σ , l , x , and $(\cdot)^\top$ represent signal standard deviation, length scale, VFZ coordinates, and transpose operator, respectively.

The values of σ and l are typically estimated from the data. One approach, employed by the `fitrgp()` routine in MATLAB involves numerical optimization via gradient descent which maximizes the likelihood as a function of these parameters [3].

An alternative approach, adapted from geostatistical applications, involves calculation of the empirical semivariogram [REF], in which one bins the space of pairwise distances between all GBs and then calculates half the average pairwise variance of the corresponding property values in each bin:

$$\kappa(d_k) = \frac{1}{2N_k} \sum_{\substack{d_\Omega(x_i, x_j) \in \\ [d_k^-, d_k^+]}} |E(x_i) - E(x_j)|^2 \quad (\text{S3})$$

where x_i and x_j are the crystallographic coordinates of GBs i and j , $d_\Omega(x_i, x_j)$ is the distance between them, and $E(x_i)$ and $E(x_j)$ are their respective energies. d_k is the location (distance) of the k -th bin center having left- and right-hand limits d_k^- and d_k^+ , and N_k is the number of measurement pairs whose distance falls in the k -th bin. Due to limited sampling of large distances and the fact that the most informative part of the semivariogram is the region near $d_k = 0$, it is customary to limit the semivariogram to half of the maximum distance [REF]. The empirical semivariogram is then fit with an analytical model to obtain the parameters of the kernel (covariance) function taking advantage of the relationship

$$\kappa(d) = \sigma_f^2 - \mathbf{C}_{\mathbf{m}_0}(d) \quad (\text{S4})$$

where we have made explicit the stationarity of the Gaussian kernel (i.e. that it depends only the distance between two points, not on their respective locations).

Having obtained a value for the length scale kernel parameter l , one can define a correlation length for the data. If the distance between two points is equal to l the kernel function indicates that their correlation will be equal to $\rho = \exp(-1/2) \approx 0.61$. However, one might reasonably want to know the length scale over which GB properties are correlated by a different amount. In general for a specified correlation strength, ρ , the corresponding correlation length is given by

$$l'(\rho) = l \sqrt{-2 \ln \rho} \quad (\text{S5})$$

We will refer to the parameter l as *the* correlation length, but one can use Eq. (S5) to determine the length scale corresponding to any specified correlation strength.

S2.2. Global Correlation Lengths

One of the most fundamental observations comes from the shape of the empirical semivariograms. As mentioned in Section 2.3, there are a wide variety of kernels that could be candidates for modeling correlations in different systems. Different kernels are used to capture different types of correlations, and each has a characteristic signature that can be observed in the semivariogram. For example, when a system exhibits exponential-type correlations, the semivariogram manifests this in the form of an exponential convergence towards an asymptotic constant value at long-distances. Linear and power-

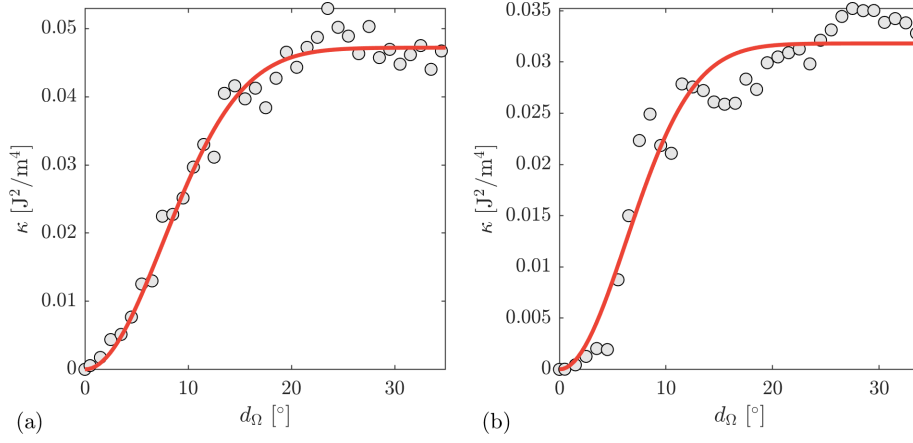


Figure S1: Empirical semivariograms (markers) for the (a) Ni and (b) Fe datasets. Solid lines show the fits of the analytical semivariogram models (Eq. S4).

type correlations manifest an absence of a long-distance plateau. In the present system, there is a clear change in concavity in the semivariogram, this is a signature of correlations that are Gaussian in nature. Thus, we find that GB energy correlations in these systems are Gaussian, and should therefore be modeled using Gaussian kernels.

S2.3. Local Correlation Lengths

The local empirical semivariograms are noisier than the global empirical semivariograms, likely due to considering fewer GB pairs. Nevertheless, reasonable fits were obtained for most of the GBs with the exception of the $\Sigma 5$ GB in the Ni dataset. In many of the local semivariograms we again see the signature change in concavity suggesting that the local correlations in the vicinity of these GBs are also Gaussian in nature. However, there are some exceptions where the nature of the correlations is more ambiguous. We anticipate that this ambiguity could be resolved with datasets that are either larger (compared to the Ni dataset) or having less noise (compared to the Fe dataset) than those considered here. However, the local empirical semivariograms seem to be generally consistent with Gaussian-type correlations.

The traditional Gaussian kernel exhibits the property of stationarity, meaning that the covariance depends only on the distance between two points, not on their respective locations (i.e. $\mathbf{C}_{\mathbf{m}_0}(x_i, x_j) = \mathbf{C}_{\mathbf{m}_0}(d(x_i, x_j))$). The use of a stationary kernel implies a prior assumption that there is a single global correlation length that applies everywhere. The fact that we observe significant variation in correlation length across the GB character space suggests that it would be better to employ non-stationary kernels (this is why, when referring to the global correlation length results presented earlier, we were careful to say that the global correlation lengths hold “on average” across the space). In particular, due to the fact that the local semivariograms do seem to be generally consistent with Gaussian-type correlations, we suggest that the non-stationary version of the Gaussian kernel [REF] may be a reasonable choice. One additional potential benefit of employing non-stationary kernels might be improved resolution of cusps in the GB energy landscape.

S3. Fe Input Dataset Characteristics

We describe the methods and results used to evaluate the quality of the noisy Fe dataset (Section S3.1) and discuss intrinsic uncertainty (Section S3.2) and overprediction bias (Section S3.3). Finally, we offer suggestions on how to improve on existing datasets (Section S3.4).

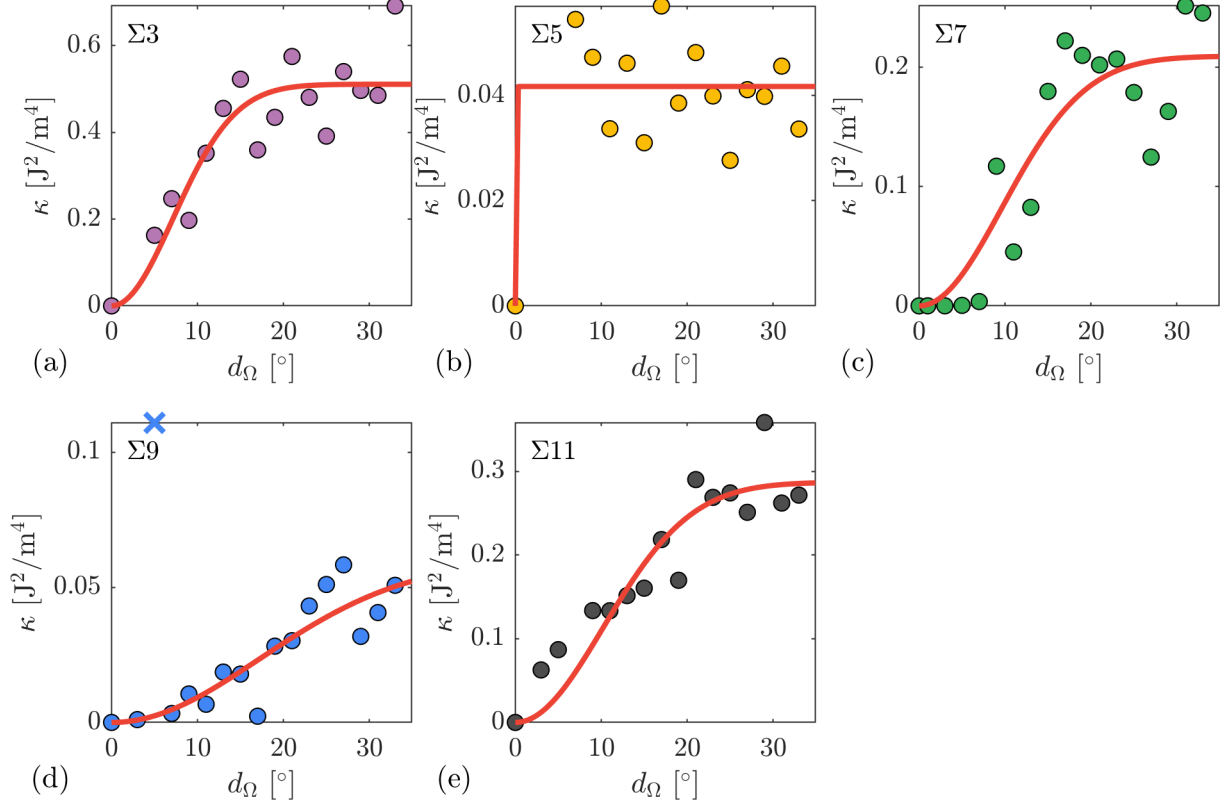


Figure S2: Local empirical semivariograms (markers), respectively centered at various low- Σ GBs for the Ni dataset. Solid lines show the fits of the analytical semivariogram models. In (d) the one point marked with an \times was considered an outlier and was excluded from the fit.

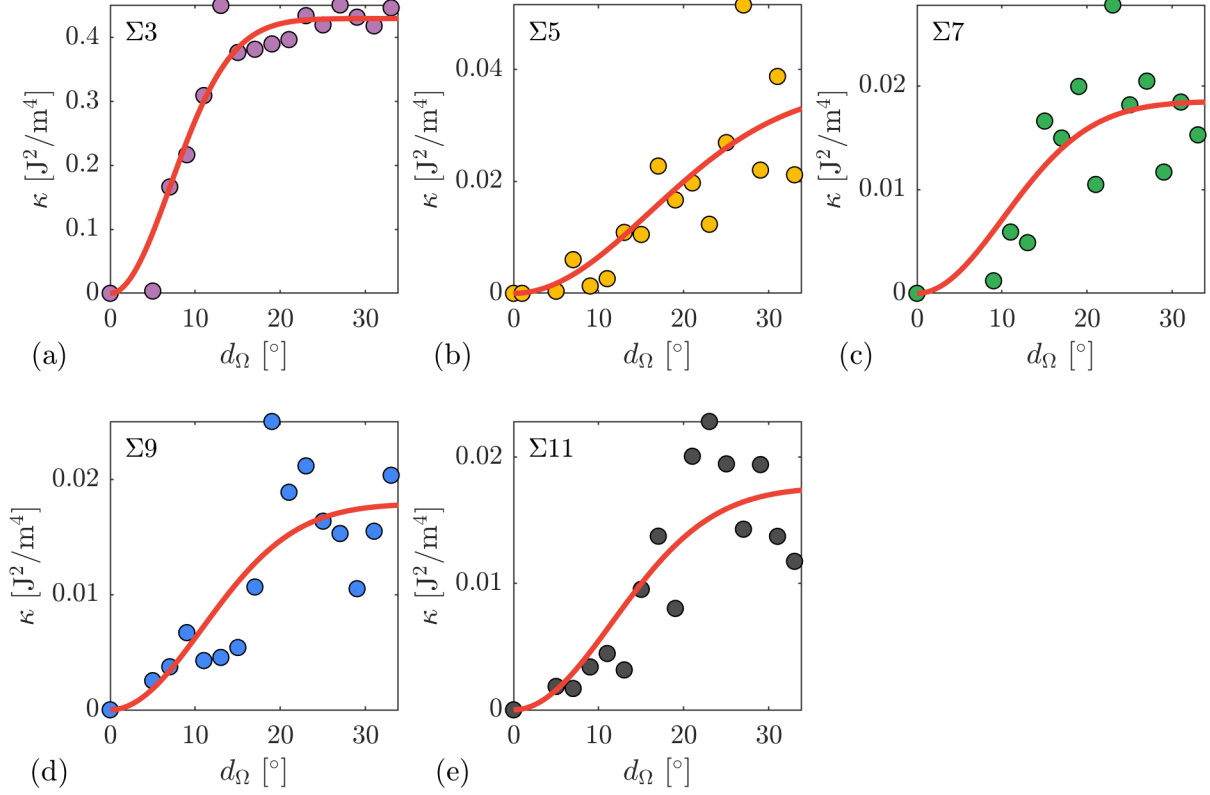


Figure S3: Local empirical semivariograms (markers), respectively centered at various low- Σ GBs for the Fe dataset. Keys for which GBs each of these correspond to in the original papers are given for Fe and Ni in [Tables S1](#) and [S2](#), respectively. Solid lines show the fits of the analytical semivariogram models.

S3.1. Methods and Results

Of the $\sim 60\,000$ GBs² in [4], $\sim 10\,000$ GBs were repeats that were identified by converting to Voronoi fundamental zone grain boundary octonions and applying VFZ repository function `avg_repeats.m`. In [4], mechanically selected GBs were those which involved sampling in equally spaced increments³ for each five degree-of-freedom parameter, and a few thousand intentionally selected GBs (i.e. special GBs) were also considered. Of mechanically and intentionally selected GBs, 9170 and 112 are repeats, respectively, with a total of 2496 degenerate sets⁴ (see Figure S4 for a degeneracy histogram). Thus, on average there is a degeneracy of approximately four per set of degenerate GBs.

By comparing GBE values of (unintentionally⁵) repeated GBs in the Fe simulation dataset [4], we can estimate the intrinsic error of the input data. For example, minimum and maximum deviations from the average value of a degenerate set are -0.2625 J m^{-2} and 0.2625 J m^{-2} , respectively, indicating that a repeated Fe GB simulation from [4] can vary by as much as 0.525 J m^{-2} , though rare. Additionally, Root mean square error and mean absolute error values can be obtained within each degenerate set by comparing against the set mean. Overall root mean square error and mean absolute error are then obtained by averaging and weighting by the number of GBs in each degenerate set. Following this procedure, we obtain an average set-wise root mean square error and mean absolute error of $0.065\,29\text{ J m}^{-2}$ and $0.061\,90\text{ J m}^{-2}$, respectively, which is an approximate measure of the intrinsic error of the data. Figure S5 shows histograms and parity plots of the intrinsic error. The overestimation of intrinsic error mentioned in the main text (Section S3.2) could stem from bias as to what type of GBs exhibit repeats based on the sampling scheme used in [4] and/or that many of the degenerate sets contain a low number of repeats (Figure S4).

Next, we see that by binning GBs into degenerate sets, most degenerate sets have a degeneracy of fewer than 5 Figure S4. We split the repeated data into sets with a degeneracy of fewer than 5 and greater than or equal to 5 and plot the errors (relative to the respective set mean) in both histogram form (Figure S5a and Figure S5c, respectively) and as hexagonally-binned parity plots (Figure S5b and Figure S5d, respectively). While heavily repeated GBs tend to give similar results, occasionally repeated GBs often have larger GBE variability. This could have physical meaning: Certain types of (e.g. high-symmetry) GBs tend to have less variation (i.e. fewer and/or more tightly distributed metastable states). However, it could also be an artifact of the simulation setup that produced this data (e.g. deterministic simulation output for certain types of GBs).

S3.2. Intrinsic Uncertainty

We estimate the intrinsic uncertainty of an Fe 0 K molecular statics simulation dataset to be $0.065\,29\text{ J m}^{-2}$ and $0.061\,90\text{ J m}^{-2}$ depending on whether root mean square error or mean absolute error estimates are used, respectively. Minimum and maximum error was -0.2625 J m^{-2} and 0.2625 J m^{-2} , respectively.

First, because only a single metastable state was used for each GBE simulation, both the training and validation data are subject to noise, consistent with a wide lateral spread of predictions and the intrinsic uncertainty estimation (Figure S5). The Fe simulation dataset Gaussian process regression mixture model gives lower root mean square error ($0.055\,035\text{ J m}^{-2}$) and mean absolute error ($0.039\,185\text{ J m}^{-2}$)

²The “no-boundary” GBs (i.e. GBs with close to 0 J m^{-2} GBE) were removed before testing for degeneracy.

³In some cases, this was equally spaced increments of the argument of a trigonometric function.

⁴A degenerate “set” is distinct from a Voronoi fundamental zone grain boundary octonion “set”, the latter of which is often used in the main text.

⁵To our knowledge, the presence of repeat GBs were not mentioned in [4] or [5]

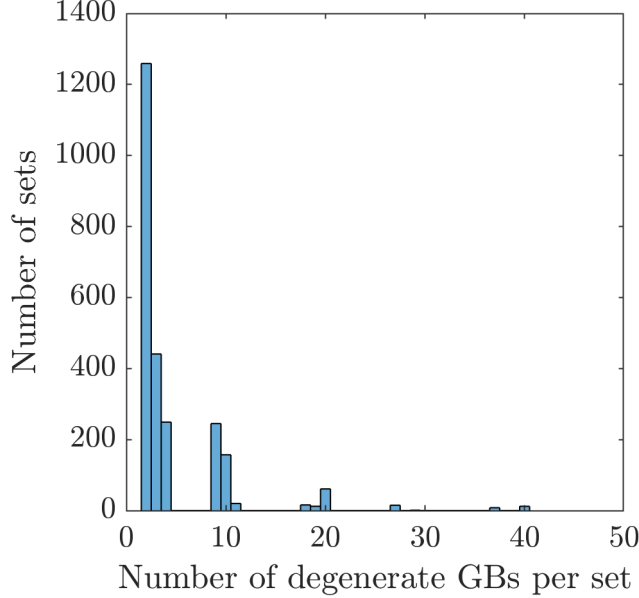


Figure S4: Histogram of number of sets vs. number of degenerate GBs per set for the Fe simulation dataset [4]. Most sets have a degeneracy of fewer than 5.

than the uncertainty estimates. This indicates that the uncertainty itself is somewhat overestimated⁶. The fact that both model and uncertainty metrics are relatively close and the prediction [1] and uncertainty parity plots (Figure S5b) are similar suggests that the model is performing well. It also suggests that further improvements in the model relative to the “true” values will be “hidden”, i.e. they will probably not manifest as lower root mean square error or mean absolute error nor as more tightly distributed parity plots despite improving performance “behind the scenes”.

S3.3. Overprediction Bias

Next, given the theoretical existence of a true minimum GBE for a given GB, the predictions which were based on metastable GBEs can be assumed to have an overprediction bias relative to the true minimum. On average, we expect this overprediction bias relative to the true minimum GBE (rather than the most likely metastable state) may be on the order of a few hundred mJ m^{-2} and may vary as a function of true minimum GBE. In other words, the model obtained is probably an estimate of the most likely metastable GBE rather than the true minimum GBE. This is akin to saying that we obtain from this data a model that approximates the non-equilibrium, Stillinger quenched red curve of Figure 4(c1) in [6], not the minimum GBE blue curve of the same chart. See [6] for an in-depth treatment of equilibrium and metastable GBE.

S3.4. Improving on Existing Datasets

Finally, datasets where multiple metastable GBEs (e.g. 3-10 repeats) are provided for each GB will likely greatly improve the performance of the Gaussian process regression model in predicting either

⁶The prediction error of a model typically cannot be less than the noise of the prediction data of a model even if the model is estimating the true prediction values with better accuracy than the noise (which is very possible and even expected with Gaussian process regression models when the noise in the input data is approximately Gaussian).

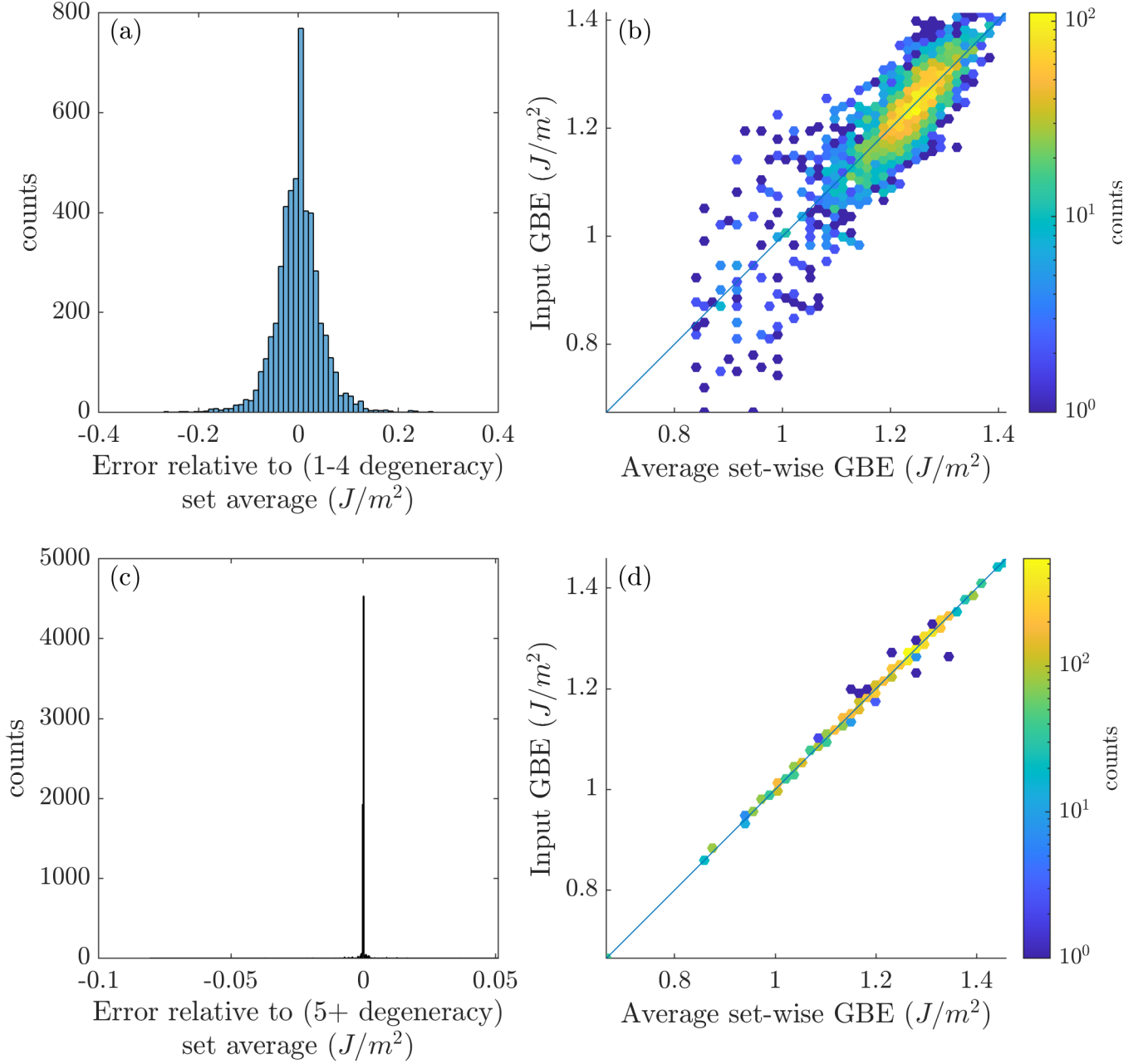


Figure S5: Degenerate GBs sets are split into those with a degeneracy of fewer than 5 and greater than or equal to 5 and plotted as (a) and (c), respectively) error histograms and (b) and (d), respectively) hexagonally-binned parity plots. Large degenerate sets tend to have very low error, whereas small degenerate sets tend to have higher error. In other words, GBs that are more likely to be repeated many times based on the sampling scheme in [4] tend to give similar results, whereas GBs that are less likely to be repeated often have larger variability in the simulation output. We do not know if this has physical meaning or is an artifact of the simulation setup.

the most likely metastable GBE (when all GBEs are considered) or the true minimum GBE (when only the minimum GBE is considered for each GB) and may even negate the need for a Gaussian process regression mixture approach. Thus, it is suggested that, where feasible, future large-scale GB bicrystal simulation studies report all property data for repeated trial runs rather than a single trial run or a single value from a set of trial runs. Ideally, data for the three additional microscopic degrees of freedom for GBs (which falls into the category of epistemic uncertainty in this work) would also be included. We believe it is likely that minimum energy paths (i.e. paths of steepest descent) in the GBE landscape depend on both macroscopic and microscopic degrees of freedom (in total, 8DOF) and could offer a more holistic view of GB behavior that better mimics and explains experimental grain growth observations. Indeed, it has been experimentally observed that at least some GB migration mechanisms involve structural transformations between equilibrium GBs via metastable states [7].

S4. Gridded Sampling for Numerical Differentiation

An isotropically sized fundamental zone may be easier to uniformly discretize than a high aspect-ratio space (i.e. a fixed discretization length can be used across all dimensions). What this doesn't describe, however, is curvature. In order to create a gridded array, which is important for numerical differentiation, a hypercube with each primary axis oriented with Euclidean dimensions is to be preferred. As curvature or misalignment is introduced as may be expected with a VFZ point cloud, GBs outside of the VFZ will necessarily be sampled; this phenomena will be exaggerated in high dimensions⁷. Fortunately, most of the information is contained in the first five dimensions after singular value decomposition transformation (Section 3.1). Thus, the latter three dimensions can likely be ignored without substantially affecting e.g. an interpolation or numerical differentiation scheme.

S5. GBs Used for Path Visualization

See Table S1 and Table S2 for the Olmsted et al. [8] and Kim et al. [4] datasets, respectively.

Table S1: Minimum Σ (Sigma) GBs and corresponding IDs used for path visualization within the original Olmsted et al. [8] dataset.

| Sigma | Olmsted ID |
|-------|------------|
| 3 | 3 |
| 5 | 169 |
| 7 | 32 |
| 9 | 21 |
| 11 | 33 |

Glossary

GB grain boundary 1, 2, 4, 6–11

⁷For perspective, a discretization into 9 segments (10 points) along each dimension will have a spacing of $\sim 7^\circ$ and require 1×10^5 grid points. In order to achieve a more reasonable grid spacing of $\sim 2^\circ$, a minimum of ~ 24 discretizations (25 points) along each dimension is necessary and will produce $\sim 1 \times 10^7$ grid points.

Table S2: Minimum Σ (Sigma) GBs and corresponding IDs used for path visualization within the original Kim et al. [4] dataset.

| Sigma | Kim Special ID |
|-------|----------------|
| 3 | 7 |
| 5 | 162 |
| 7 | 259 |
| 9 | 315 |
| 11 | 406 |

GBE grain boundary energy [2](#), [7](#), [8](#), [10](#)

GBO grain boundary octonion [1](#), [2](#)

VFZ Voronoi fundamental zone [1-3](#), [7](#), [10](#)

References

- [1] S. G. Baird, E. R. Homer, D. T. Fullwood, O. K. Johnson, Five Degree-of-Freedom Property Interpolation of Arbitrary Grain Boundaries via Voronoi fundamental zone Framework, Computational Materials Science (Under Review) 26.
- [2] T. Francis, I. Chesser, S. Singh, E. A. Holm, M. De Graef, A geodesic octonion metric for grain boundaries, Acta Materialia 166 (2019) 135–147. doi:[10.1016/j.actamat.2018.12.034](https://doi.org/10.1016/j.actamat.2018.12.034).
- [3] Exact GPR Method - MATLAB & Simulink, <https://www.mathworks.com/help/stats/exact-gpr-method.html>, ????
- [4] H.-K. Kim, S. G. Kim, W. Dong, I. Steinbach, B.-J. Lee, Phase-field modeling for 3D grain growth based on a grain boundary energy database, Modelling Simul. Mater. Sci. Eng. 22 (2014) 034004. doi:[10.1088/0965-0393/22/3/034004](https://doi.org/10.1088/0965-0393/22/3/034004).
- [5] H. K. Kim, W. S. Ko, H. J. Lee, S. G. Kim, B. J. Lee, An identification scheme of grain boundaries and construction of a grain boundary energy database, Scripta Materialia 64 (2011) 1152–1155. doi:[10.1016/j.scriptamat.2011.03.020](https://doi.org/10.1016/j.scriptamat.2011.03.020).
- [6] J. Han, V. Vitek, D. J. Srolovitz, Grain-boundary metastability and its statistical properties, Acta Materialia 104 (2016) 259–273. doi:[10.1016/j.actamat.2015.11.035](https://doi.org/10.1016/j.actamat.2015.11.035).
- [7] J. Wei, B. Feng, R. Ishikawa, T. Yokoi, K. Matsunaga, N. Shibata, Y. Ikuhara, Direct imaging of atomistic grain boundary migration, Nat. Mater. (2021). doi:[10.1038/s41563-020-00879-z](https://doi.org/10.1038/s41563-020-00879-z).
- [8] D. L. Olmsted, E. A. Holm, S. M. Foiles, Survey of computed grain boundary properties in face-centered cubic metals-II: Grain boundary mobility, Acta Materialia 57 (2009) 3704–3713. doi:[10.1016/j.actamat.2009.04.015](https://doi.org/10.1016/j.actamat.2009.04.015).