# Scientific Benchmarks for Guiding Macromolecular Energy Function Improvement

Andrew Leaver-Fay*, *Matthew J. O'Meara**, Mike Tyka, Ron Jacak, Yifan Song, Elizabeth H. Kellogg, James Thompson, Ian W. Davis, Roland A. Pache, Sergey Lyskov, Jeffrey J. Gray, Tanja Kortemme,, Jane S. Richardson, James J. Havranek, Jack Snoeyink, David Baker, Brian Kuhlman

# Score Function Consensus

**New Terms**

- dun10
- orbitals
- H-patch
- cart_bonded/mm_*
- pH
- hackelec
- gb/pb
- geometric solvation
- env_dep reference weights

**Re-parametrizations**

- score12'
- -correct
- sp2 hbond
- lennard jones cutoff/radii
- softrep/hardrep
- P(aa|pp)
- pair

# How to Demonstrate Improvement

- Explain problem/solution
- Run scientific benchmarks

# Outline

- Two tools
  - Features Analysis
  - OptE
- Example Modification
  - Score12bicubic
  - Dun10
- Scientific Benchmarks

# Boltzmann Distribution
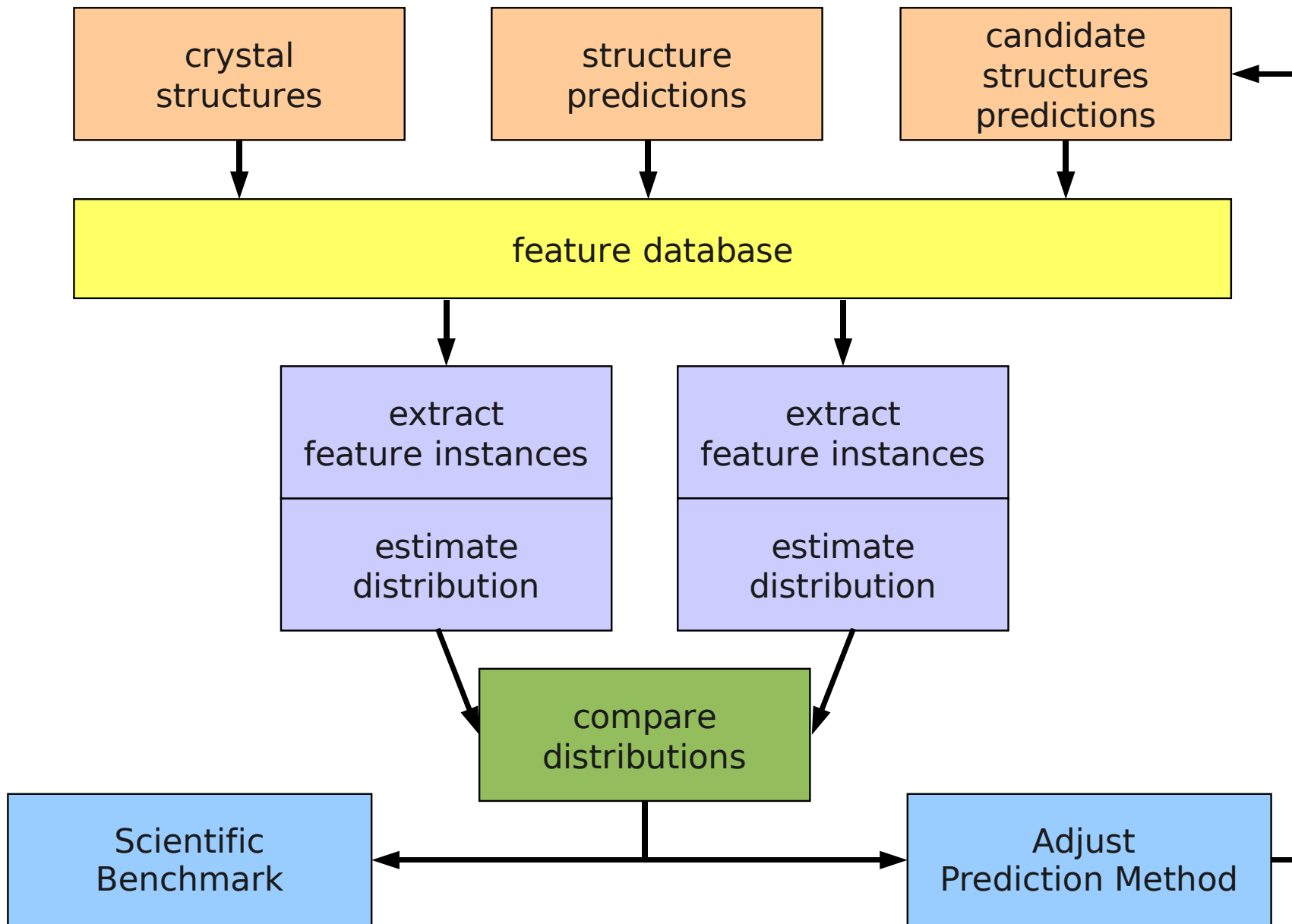
$$P(C) = \frac{1}{Z} e^{\frac{-E(C)}{kT}}$$
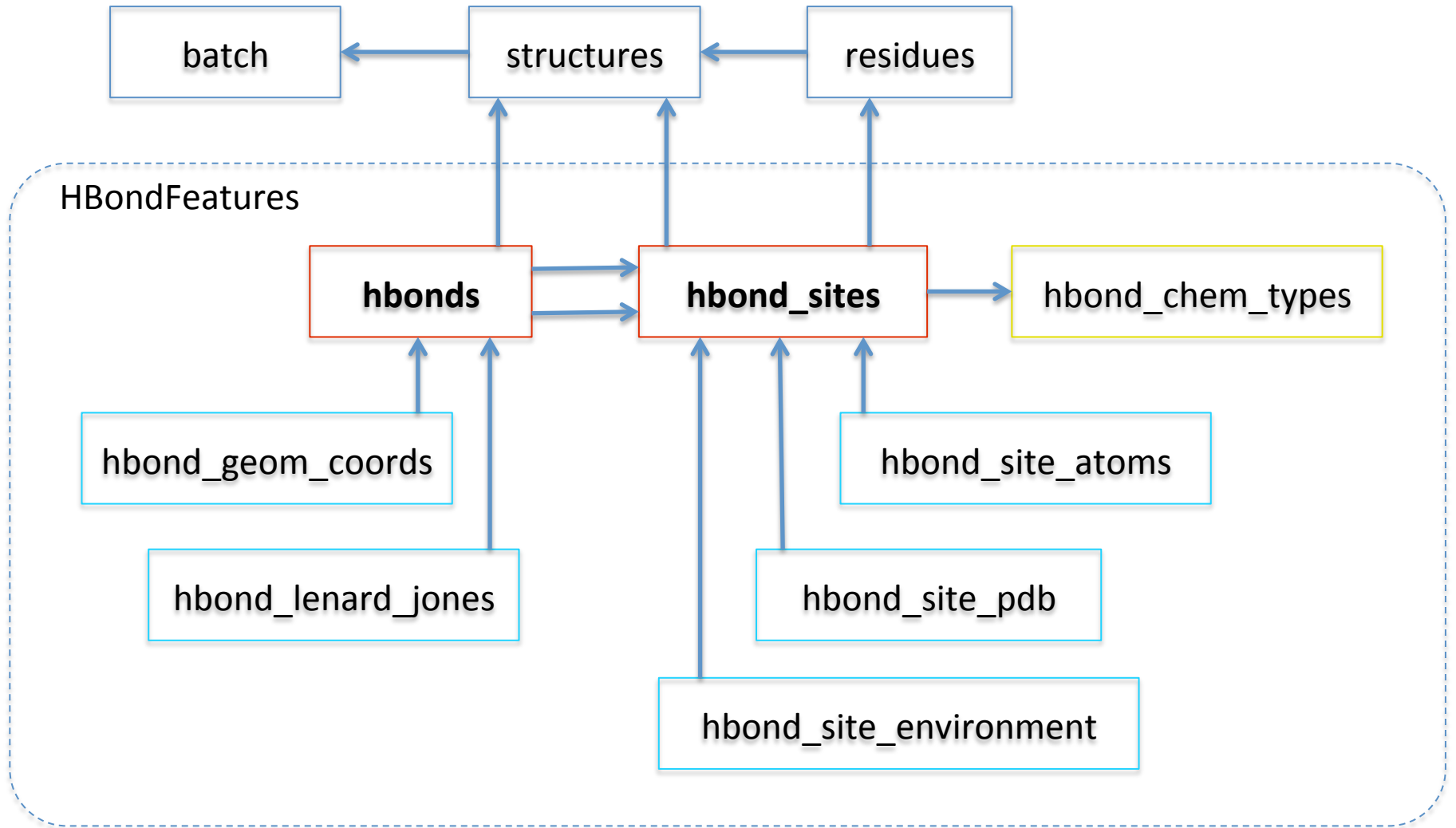
## What should P(C) be?

# Feature Analysis

**Feature**: a geometric observable of a molecular conformation

A distribution over *conformation* space ⟶ A distribution over *feature* space
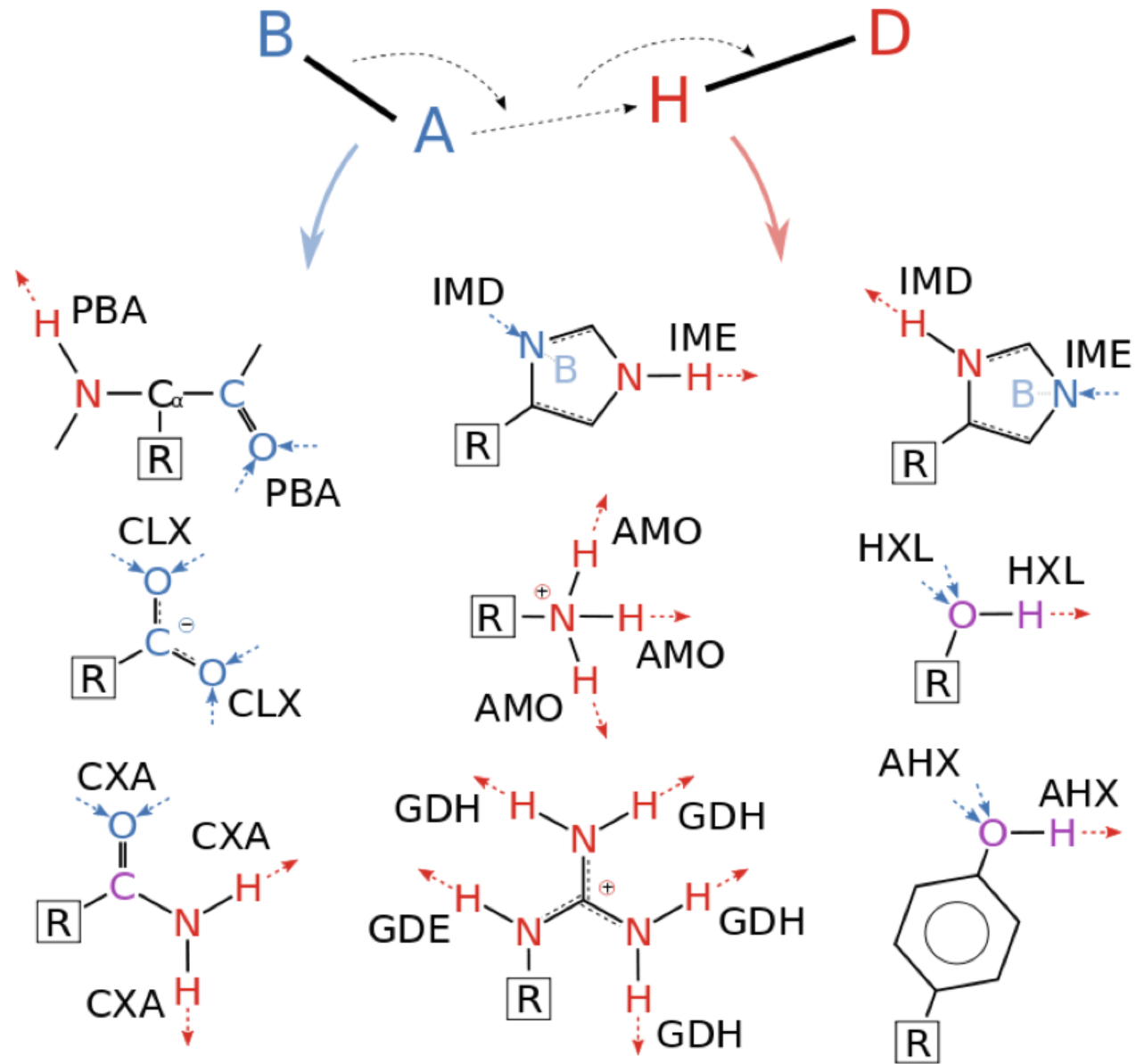
# Feature-Based Workflows

# Hbond Features Reporter
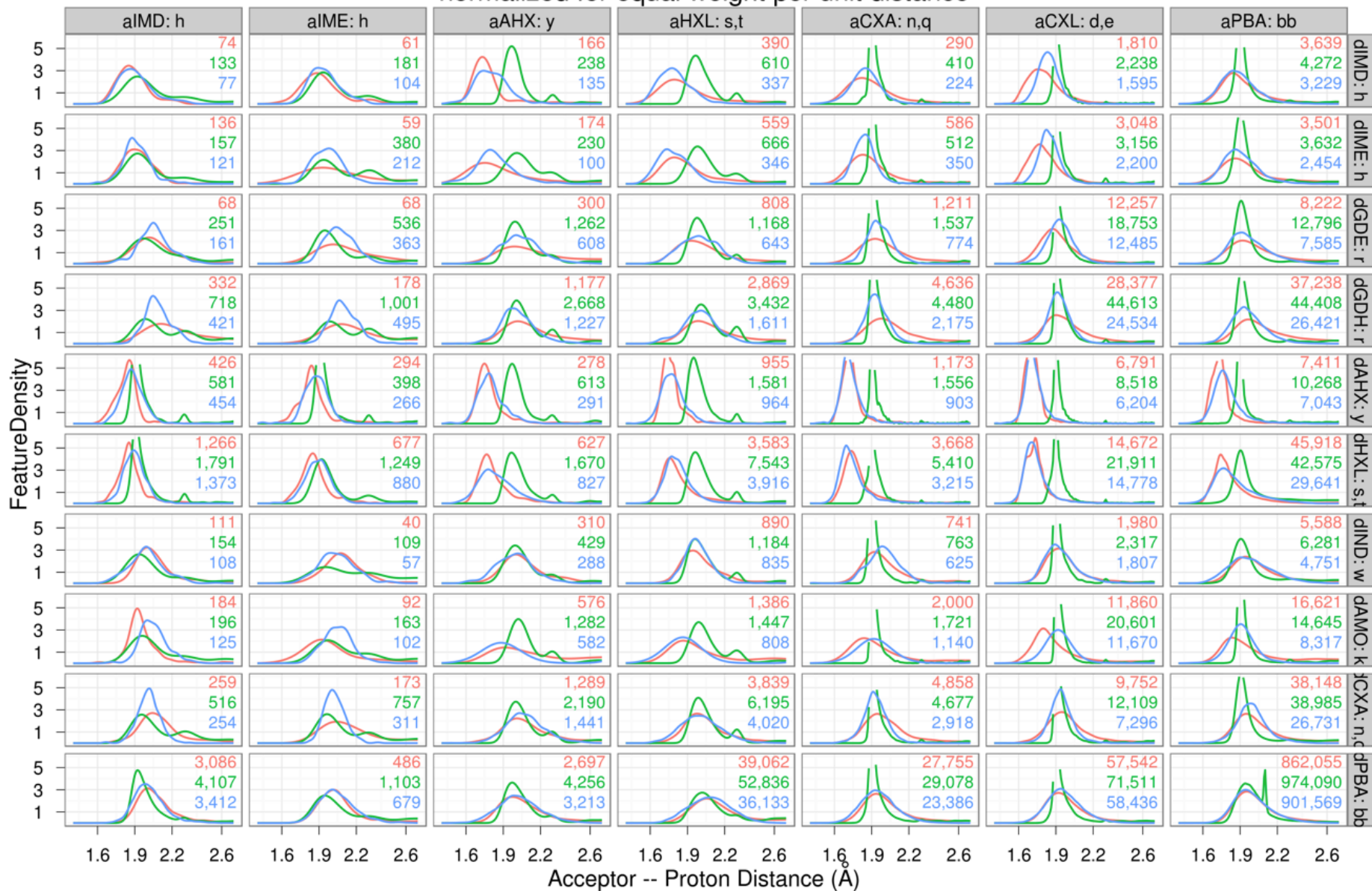
# Feature Reporters => "Schema Parts"

| Meta | One Body | Two Body | Multi Body |
|------|----------|----------|------------|
| Protocol | Residue | Pair | Structure |
| Batch | ResidueConformation | AtomAtomPair | PoseConformation |
| JobData | ProteinResidueConformation | AtomInResidue- | RadiusOfGyration |
| PoseComments | ProteinBackboneTorsionAngle | AtomInResiduePair | SecondaryStructure |
| | ResidueBurial | ProteinBackbone- | HydrophbicPatch |
| **Experimental Data** | ResidueSecondaryStructure | AtomAtomPair | Cavity |
| PdbData | GeometricSolvation | HBond | GraphMotif |
| PdbHeaderData | BetaTurn | Orbital | SequenceMotif |
| DDG | RotamerBoltzmannWeight | SaltBridge | Rigidity |
| NMR | ResidueStrideSecondaryStructure | LoopAnchor | VoronoiPacking |
| DensityMap | HelixCapping | DFIREPair | InterfaceAnalysis |
| MultiSequenceAlignment | BondGeometry | ChargeCharge | |
| HomologyAlignment | ResidueLazaridisKarplusSolvation | | **Energy Function** |
| | ResidueGeneralizedBornSolvation | **Multi Structure** | ScoreFunction |
| **Chemical** | ResiduePoissonBoltzmannSolvation | ProteinRMSD | ScoreType |
| AtomType | Pka | ResidueRecovery | StructureScores |
| ResidueType | ResidueCentroids | ResiduePairRecovery | ResidueScores |
| | | ResidueClusterRecovery | HBondParameters |
| | | Cluster | <EnergyTerm>Parameters |

# Refined Hydrogen Bond Model
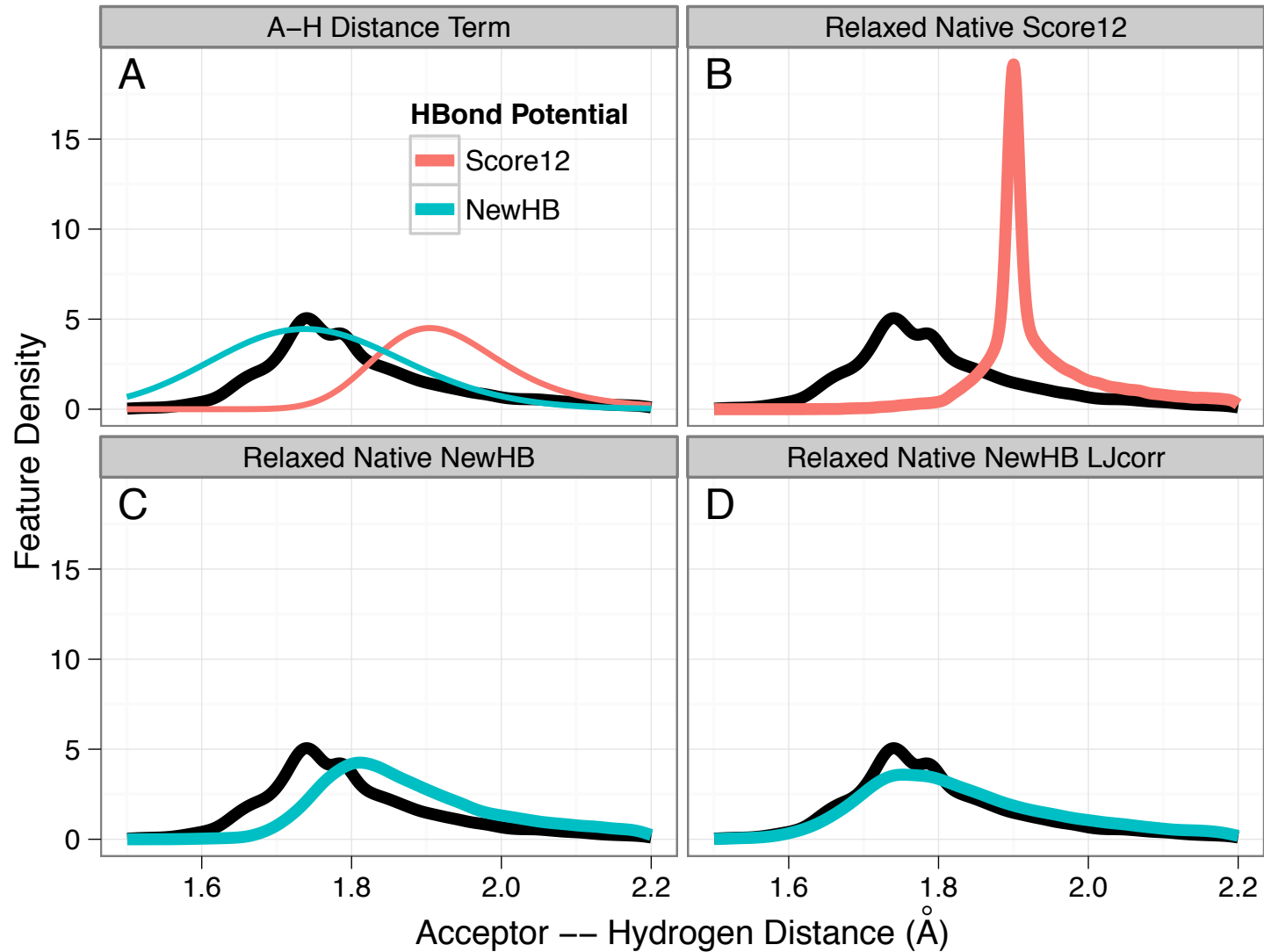
HBond A-H Distance by Chemical Type, B-Factor < 30
normalized for equal weight per unit distance

FeatureDensity

Acceptor -- Proton Distance (Å)

sample_source

top8000_r46440_111212    top8000_relax_r46440_111213    top8000_relax_olf_r46015_111119
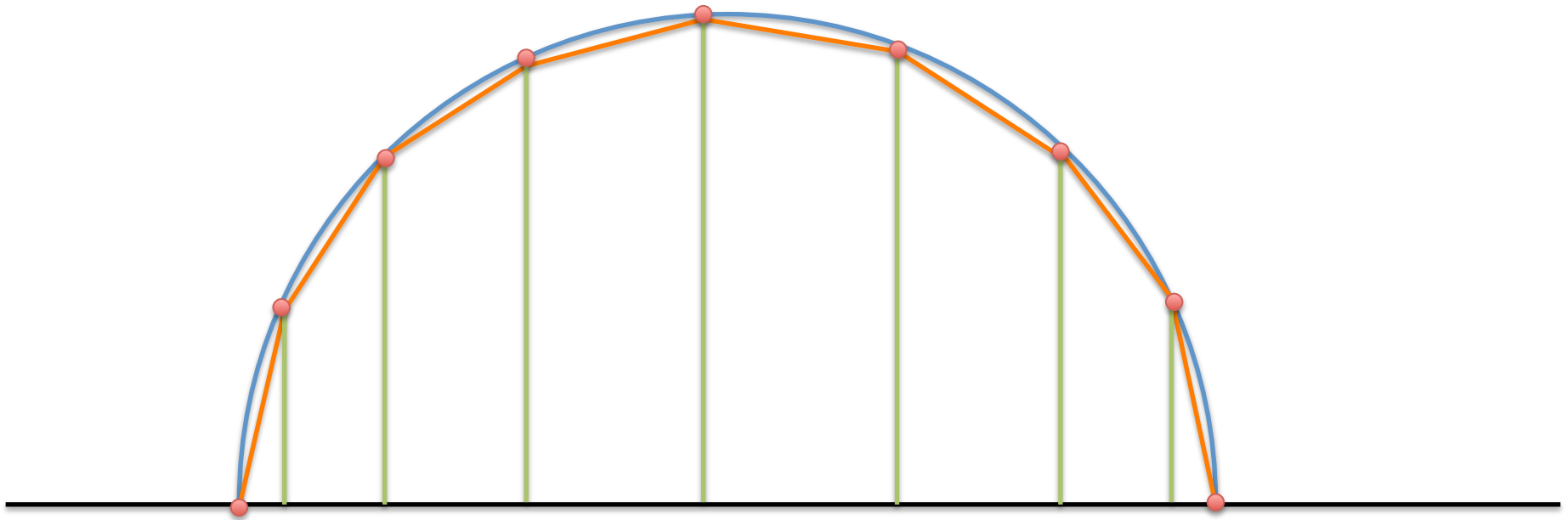
# HBond Potential

# OptE Overview

- Encode Scientific Benchmarks as a "Loss Function"
  - Seq.-Profile Recovery, Rot. Recovery, ddG, etc.

- Optimize weights to minimize loss
  - Estimate Partition Function
  - Optimize Weights
  - Repeat

- Jim, Andrew, Ron, Liz, Yifan, Mike, James, Ian, more...
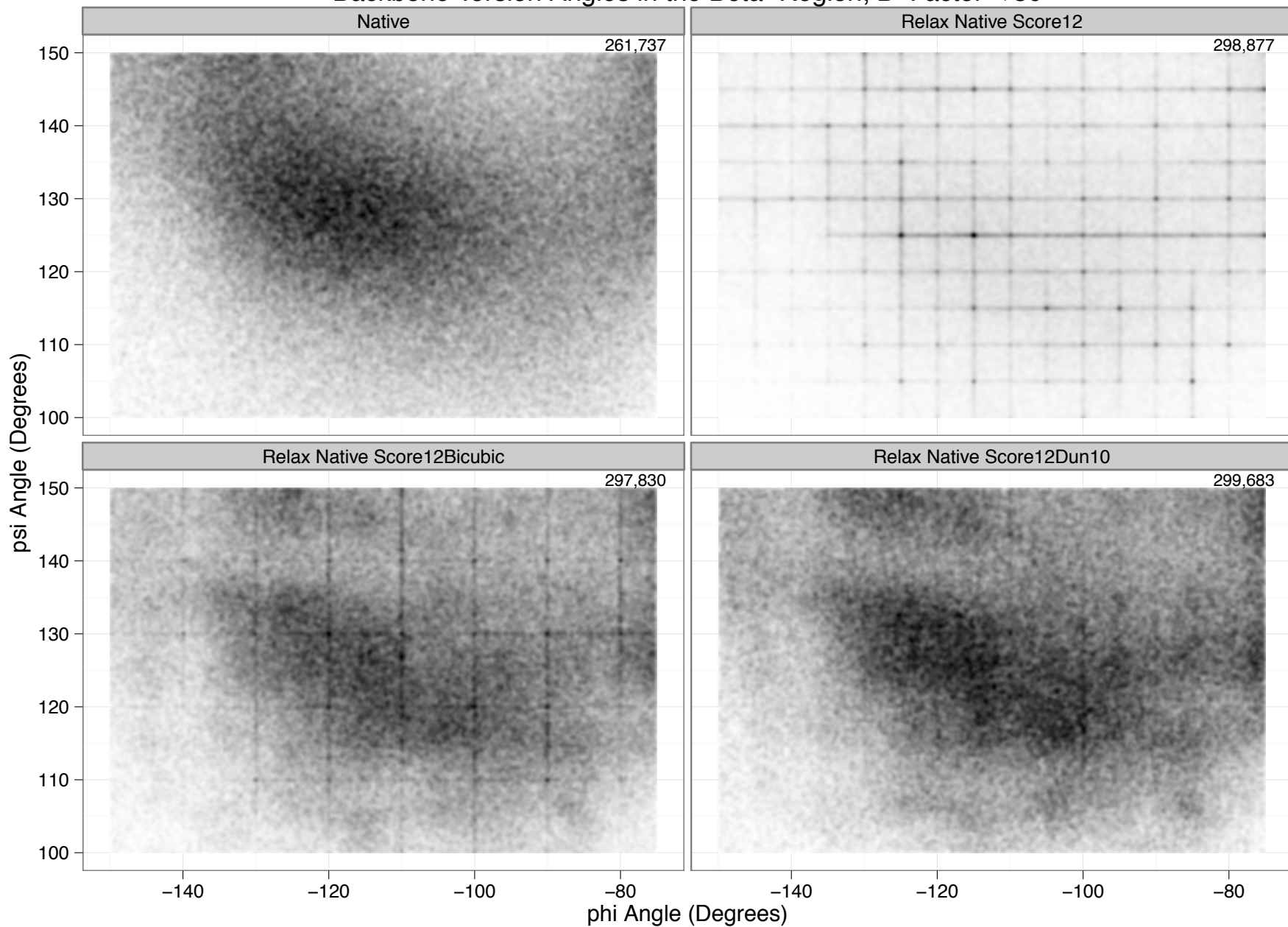
# OptE Capabilities

- Re-fit reference weights
  - Sequence-profile-recovery

- Test targeted hypotheses with a few weights
  - Hackeleck
  - CH-bond potential
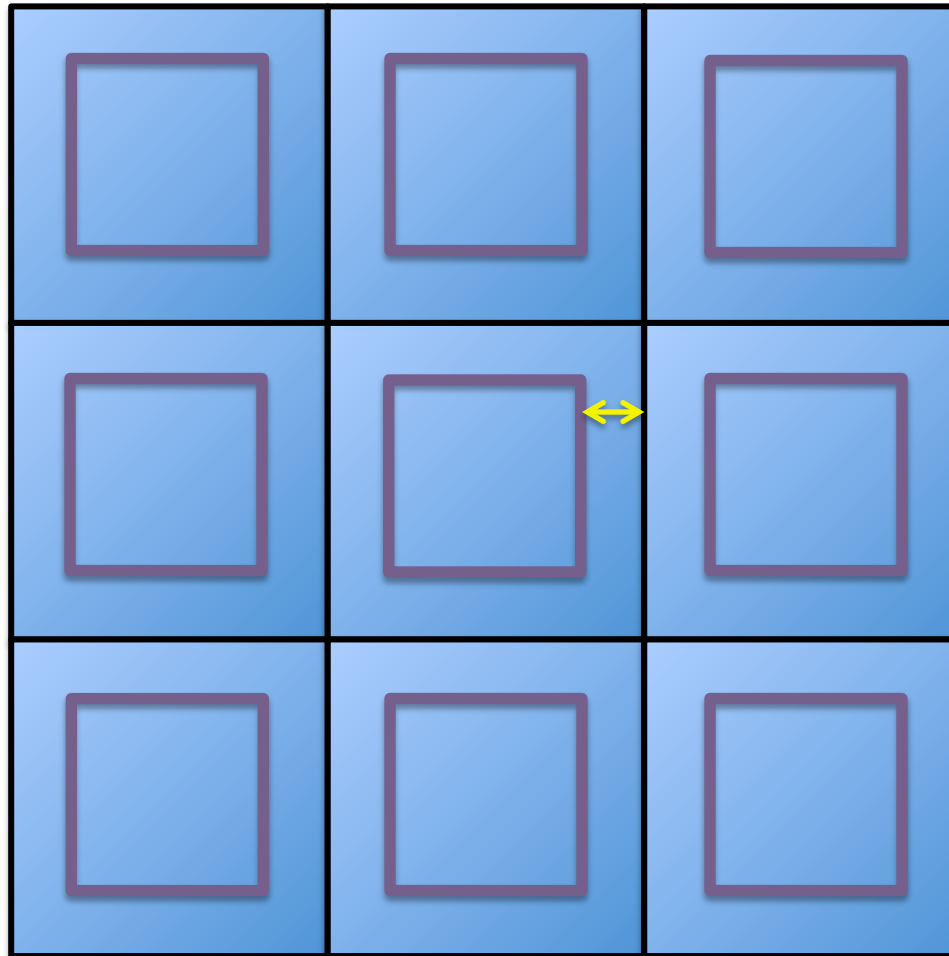
# Linear vs Spline Interpolation

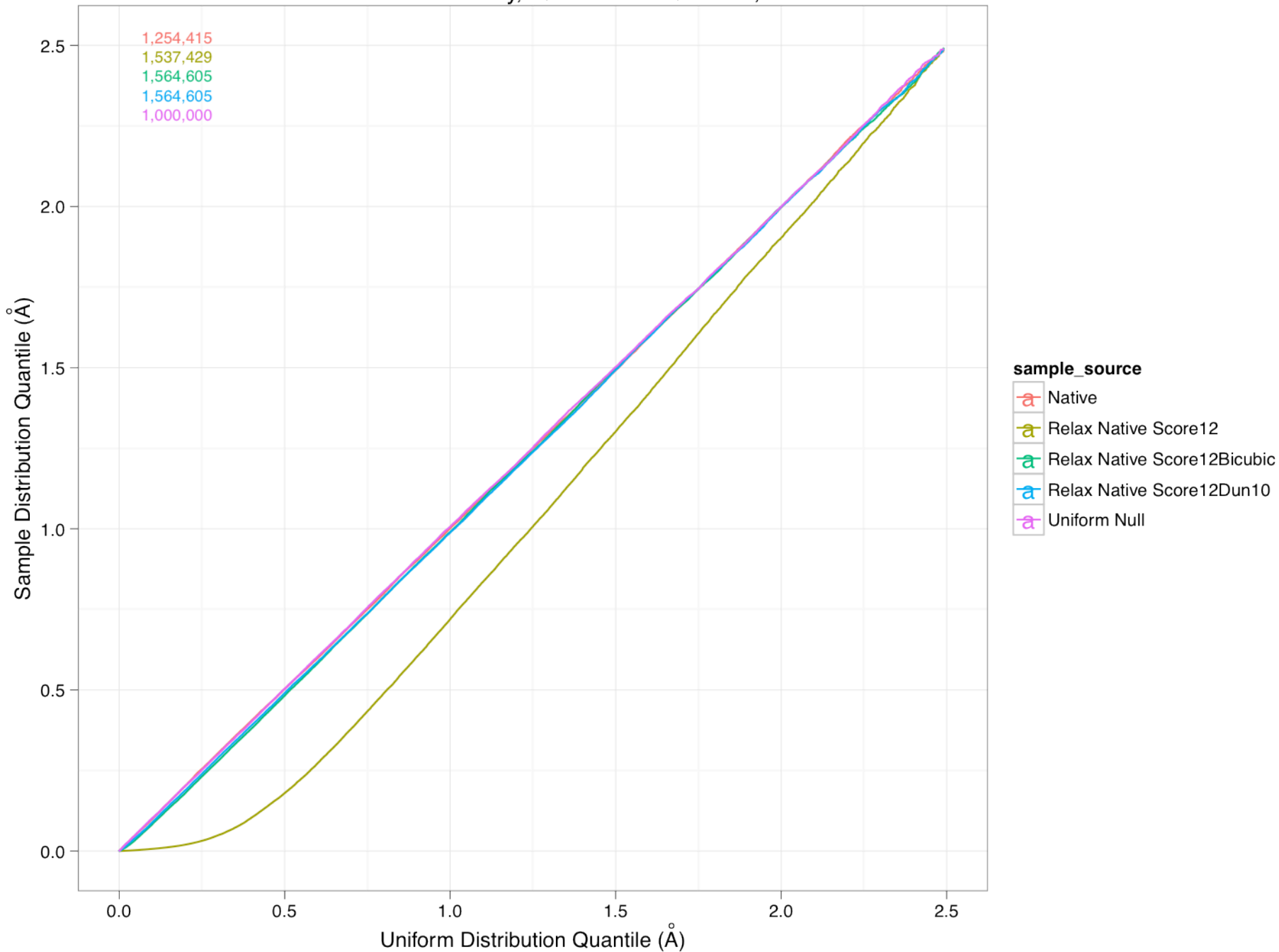Backbone Torsion Angles in the Beta−Region; B−Factor < 30

# off_grid feature

Min Distance From Grid Boundary, Quantile vs Quantile; B−Factor < 30

# Dun10: Semi Rotameric
# $P(X_2 \mid X_1 = \text{trans})$

# Major Scientific Benchmarks

**Currently Available**

- Rotamer Recovery
- Sequence Recovery
- ddG Prediction
- Loop Recovery
- Ab-relex Recovery
- Docking local-refine
- RNA benchmark(s)

**Upcoming**

- Fit into electron density (frank)
- LoopHash discrimination (TJ)
- Ligand docking (Rocco/Sagar/Ora)
- Single Mutant Scan (Yifan)
- Fix-interface Design (Jacob)
- Flex-Interface Design (Sarel)
- Flex-BB Design (Nobu)
- NMR recovery (Oliver)
- <Your Benchmark here>

score12 : score vs loopcarms  by  target
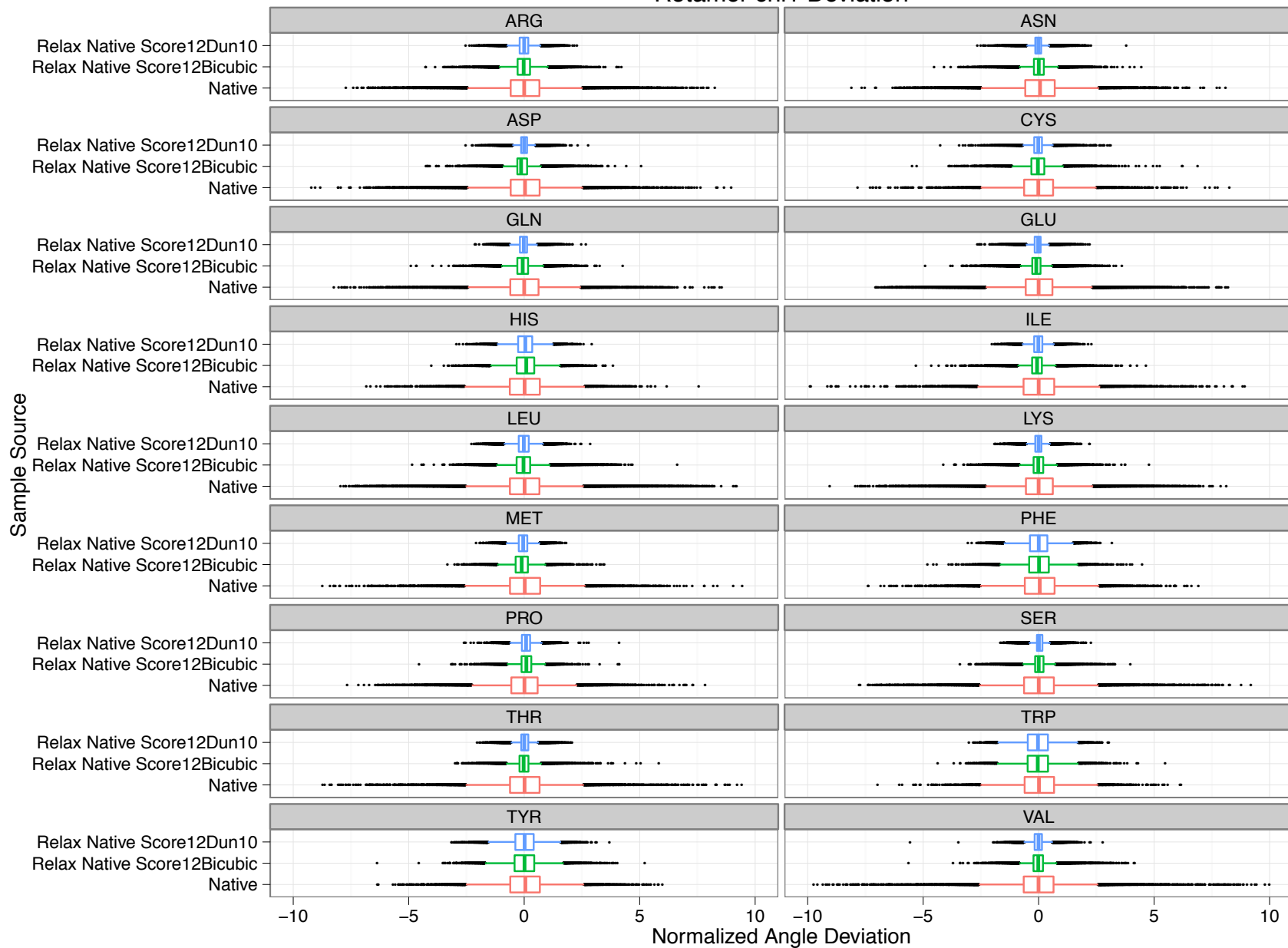
# Bicubic/Dun10 Results

| Energy Function | Rotamer Recovery Benchmark | | | | Seq. Rec. Bench | | ΔΔG Bench | High-Res. Refinement Benchmark | | | Loop Modeling Benchmark | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | pack rots (%) | min pack (%) | rot. trials (%) | rt-min (%) | % Rec | KL-Div. | R-Value | #(pNat > 0.8) | Σ pNat | #(eNat < eDec ) | 1st Quart. (Å) | Med (Å) | 3rd Quart. (Å) |
| Score12 | 66.19 | 69.07 | 71.49 | 73.12 | 32.6 | 0.019 | 0.69 | 67 | 74.6 | 104 | 0.468 | 0.637 | 1.839 |
| Score12' | - | - | - | - | 37.0 | 0.008 | 0.67 | - | - | | - | - | - |
| Score12Bicubic | 66.24 | 67.51 | 71.52 | 73.15 | 37.6 | 0.010 | 0.68 | 68 | 77.9 | 105 | 0.499 | 0.644 | 1.636 |
| Score12Dun10 | 67.82 | 70.50 | 72.60 | 74.23 | 37.6 | 0.009 | 0.67 | 60 | 72.0 | 104 | 0.461 | 0.677 | 1.463 |

# Thanks

- Brian Kuhlman / Jack Snoeyink
- Andrew Leaver-Fay
- Rosetta Community

Rotamer chi1 Deviation

# Top8000 Data Set

- Assembled by Richardson Lab (at Duke)
  - March 2011 snapshot of **Protein Databank**
  - Clustered so intra-cluster homology is at most 70%
  - Filter out structures having
    - greater than 2A resolution
    - known oddities
    - non-canonical amino acids
  - Selected best average MolProbity score and resolution
  - Place hydrogen atoms using Reduce
  - 6,563 Chains

# Rotamer Recovery

- ## 152 chains, 17,463 residues
  - Subset of  Top5200 (from Richardson Lab)
  - 50-200 residues each
  - at most 70% seq. homology
  - at most 1.2A resolution
- ## Recovered: all angles within 20° of native
- ## Starting info / DOFs:
  - RotamerTrials          PackRotamers
  - RTMin                     MinPack

# Sequence Recovery

- 38 large structures (Ding & Dekholyan 2006)
- Accuracy:
  - Native AA
  - Kullback-Leibler divergence with input AA-profile
- Starting info / DOFs:
  - Fixed backbone
  - PackRotamers protocol

# High Resolution Refinement

- 114 Sequences (Tyka, *et al.* 2010)
  - 4 centroid mode data sets
    - homolog fragments / relax to low RMSD + low energy
    - 6,000 FastRelax predictions
- Accuracy:
  - Boltzmann weighted probability of "near-native"
    - where less than 2A RMSD => near-native
  - If 80% near-native
  - If min near-native energy < min decoy energy

# Loop Prediction

- 45 12-Residue Loops (Mandell *et al*. 2009)
  - 8,000 Kinimatic Loop Closure (KIC) predictions
- Accuracy:
  - Min Cα-RMSD over 5 lowest-energy structures

# ddG Prediction

- 1210 Point Mutants (Kellogg *et al*. 2011)
  - Native crystal structures
  - Experimental ddG of folding
- Accuracy:
  - correlation coefficient
- Input info / DOFs:
  - all-atom, soft-rep repacking
  - backbone + sidechain, hard-rep minimization
  - uniform constraints