
IOBP2-Bench: A Public AI/ML-Ready Benchmark from a Pivotal-Class Adaptive Automated Insulin Delivery Trial

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Posing several of the most consequential ML problem formulations for adaptive
2 Automated Insulin Delivery (AID) — offline imitation of a deployed controller,
3 controller-policy sensitivity analysis under realistic input distributions, and pre-trial
4 responder prediction — requires a public dataset that simultaneously satisfies four
5 conditions: (i) a controller that *adapts during the trial* rather than a fixed-gain
6 PID/MPC; (ii) the controller’s *decision-tick log released at native cadence*; (iii) a
7 *cohort large enough for subject-level train/val/test splits with stratification*; and (iv)
8 *longitudinal coverage* long enough that adaptation dynamics are observable. To
9 our knowledge, no prior public AID release satisfies all four. IOBP2-BENCH is the
10 first benchmark that does. Built from the public Jaeb release of the Insulin-Only
11 Bionic Pancreas (IOBP2) Pivotal Trial, it covers 440 randomized participants over
12 91 days (11.5 M aligned 5-minute bins; 13 channels), ships five reference tasks
13 with at least three baselines each, and is audited end-to-end against the published
14 trial: arm-level Time-in-Range from our pipeline matches the headline NEJM
15 endpoint within ~ 1.5 percentage points (Control TIR 54.3% vs. $\sim 54\%$ in Russell
16 et al. [12]). Three findings emerge from the baselines themselves. First, the first
17 14 days of CGM are sufficient to forecast 13-week glycemic outcome (test-set
18 AUROC 0.83, 95% CI 0.71–0.92), supporting short-window responder prediction
19 for trial design. Second, on offline imitation of the iLet’s adaptive controller, basal
20 adaptation is statistically learnable from history alone whereas meal anticipation is
21 not without the announcement channel. Third, the qualitative meal-size annotation
22 carries near-zero standalone signal — the system’s tolerance to qualitative reporting
23 arises from controller adaptation, not label informativeness. Code, splits, datasheet,
24 and reproducibility scripts are released openly; processed Parquet artifacts follow
25 a credentialed-access pattern (akin to MIMIC-III/IV), reproducible from the raw
26 Jaeb release in ~ 10 min on a single CPU. All baseline numbers are reported with
27 subject-level bootstrap 95% confidence intervals.

28 1 Introduction

29 Automated Insulin Delivery (AID) systems — continuous glucose monitor (CGM) plus insulin pump
30 plus control algorithm — now form the standard of care for type-1 diabetes (T1D) in many high-
31 income settings, with at least three commercial systems shipping in the United States and growing
32 international uptake [3, 12]. Despite this clinical maturity, the public dataset landscape that ML
33 researchers can build on for AID-relevant tasks remains thin. OHIO T1DM, the de facto benchmark
34 for short-horizon CGM forecasting, contains twelve participants in its combined 2018+2020 release
35 [8, 9]; the SHANGHAI T1DM release provides 12 [14]; the recently-released T1DEXI extends to

36 several hundred participants but is centred on a six-week structured-exercise protocol rather than
37 ongoing closed-loop therapy [11].

38 **The capability gap.** Several core ML problem formulations for adaptive AID — offline imitation
39 of a deployed controller, sensitivity analysis of controller policies under realistic input distributions,
40 and pre-trial responder prediction — share dataset prerequisites that no prior public release satisfies
41 simultaneously: (i) a controller that *adapts* during the trial (not a fixed-gain PID/MPC); (ii) its
42 *decision-tick log* released at native cadence (not coarse delivery summaries); (iii) a cohort large
43 enough for subject-level train/val/test splits; and (iv) longitudinal coverage long enough to expose
44 adaptation dynamics. OHIO T1DM and SHANGHAI T1DM fail (i)–(iii); T1DEXI fails (i)–(ii); the
45 Jaeb DCLP3 release uses a fixed-gain hybrid PID-MPC with only limited decision logging. As
46 a consequence, core AID-policy ML research has had to fall back on simulator data (typically
47 UVA/PADOVA [5], single-meal-validated and unable to capture inter-day adaptation), with a known
48 sim-to-real gap.

49 The Insulin-Only Bionic Pancreas (IOBP2) Pivotal Trial [12] is the first AID trial whose public
50 release satisfies all four conditions. The iLet system uses *no* carbohydrate counting, basal-rate
51 program, or insulin-to-carb ratio: initialization requires only the participant’s body weight, meals
52 are announced qualitatively (“less / typical / more” against a learned per-period reference), and
53 the controller adapts continuously to each patient. The trial’s release through the Jaeb Center for
54 Health Research exposes the *full minute-resolution decision log of the deployed adaptive controller*
55 alongside CGM and patient-reported outcomes, for all ~ 440 randomized participants over 13 weeks.
56 Using this data for ML research, however, requires substantial engineering: 60 pipe-delimited text
57 tables across multiple sampling cadences, with mandatory attribution clauses and several non-obvious
58 encoding decisions (e.g., the qualitative meal label is hidden inside a numeric composite field whose
59 semantics are documented only in a separate glossary).

60 IOBP2-BENCH closes this gap. We provide:

- 61 1. A reproducible, fully-tested pipeline that converts the raw Jaeb release into per-subject Parquet
62 trajectories on a 5-minute uniform grid, with explicit mask channels, decoded meal-event compos-
63 ites, dropped quasi-identifiers, and protocol-faithful CGM metrics computed per Battelino et al.
64 [1] and Bergenstal et al. [2];
- 65 2. Deterministic 60/20/20 subject-level splits stratified by treatment arm, age stratum, and baseline
66 HbA1c tertile, shipped as a JSON file so that across-paper comparisons are exact;
- 67 3. Five reference tasks (Section 5) ranging from short-horizon CGM forecasting (T1) and hy-
68 poglycemia warning (T2), through postprandial incremental-AUC prediction (T3) and offline
69 imitation of the iLet’s continuously-adapting controller (T4), to 13-week responder prediction
70 (T5), each accompanied by three or more baselines and standardized evaluation;
- 71 4. An end-to-end validation against the published trial: the arm-level CGM summaries that fall out of
72 our pipeline match the headline numbers of Russell et al. [12] within sampling noise (Section 7),
73 providing a strong signal that the upstream alignment and decoding are correct;
- 74 5. A complete Datasheet for Datasets [6] with the trial’s mandatory attribution clause, an explicit
75 quasi-identifier-removal step, and a documentation of the non-obvious encoding decisions required
76 to use the raw release.

77 Three findings emerge from the baselines themselves (Section 8): (1) the first 14 days of CGM
78 forecast 13-week glycemic outcome at AUROC 0.83 (95% CI 0.71–0.92); (2) basal-insulin adaptation
79 is statistically learnable from history alone whereas meal-bolus anticipation is not without the
80 announcement channel; and (3) the qualitative meal-size annotation has near-zero standalone signal
81 — the trial’s tolerance to qualitative meal reporting is attributable to controller adaptation, not label
82 informativeness.

83 **Release plan and access model.** The Jaeb release is openly downloadable under a click-through
84 data-use agreement and a mandatory attribution clause. We adopt a deliberate *credentialed-access*
85 release pattern, comparable to the access model used by MIMIC-III/IV for de-identified ICU data: (i)
86 all conversion scripts, splits, datasheet, and baseline code are released openly under an MIT-style
87 license; (ii) users download the upstream Jaeb release directly from the Jaeb portal (preserving the
88 trial sponsor’s audit trail and the mandatory attribution clause), and (iii) re-running our pipeline
89 reproduces the 5-min-aligned Parquet trajectories bit-for-bit in approximately 10 minutes on a single

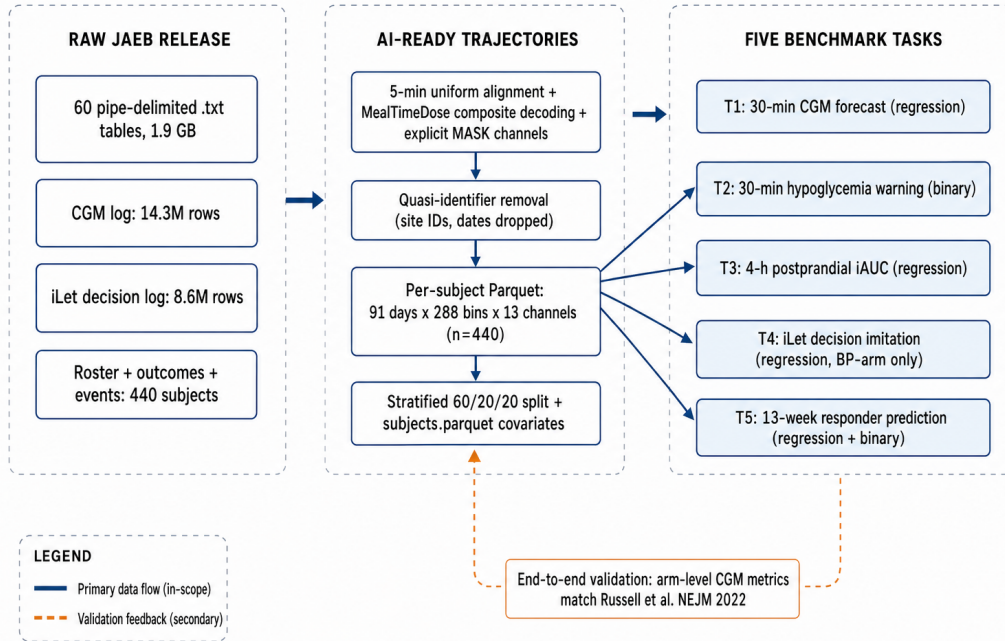


Figure 1: The IOBP2-BENCH construction pipeline. The raw Jaeb release (60 tables) is converted into per-subject 5-min-aligned Parquet trajectories with explicit masks, decoded meal events, and quasi-identifiers stripped; five reference tasks are derived from those trajectories. The dashed feedback loop (end-to-end audit): our arm-level CGM summaries match Russell et al. [12].

90 CPU. Re-distribution of the processed Parquet artifacts as a stand-alone download (i.e., bypassing
 91 the Jaeb portal) is pending written sponsor confirmation; if granted, the artifacts will be published
 92 on Zenodo with a citable DOI. Until then, the canonical access path is portal-plus-scripts, and the
 93 included `Makefile` reduces user effort to a single command.

94 **Scope and what this paper is not.** IOBP2-BENCH is a *secondary-use, ML-ready repackaging* of
 95 an already-public clinical trial release, not new clinical evidence. We do not claim novel medical
 96 findings; the headline trial results have already been reported in the *New England Journal of Medicine*
 97 [12]. IOBP2-BENCH is also *not* a closed-loop replay environment: we report open-loop imitation
 98 MAE for T4. A full simulator-coupled evaluation harness using the FDA-accepted UVA/PADOVA
 99 T1D simulator [5] is in scope for v0.2 of the benchmark and is discussed in Section 9.

100 2 Related Work

101 **Public CGM and AID datasets.** Section 1 stated the four Gap-Chain conditions and which prior
 102 public releases fail which ones. Table 1 makes the comparison explicit across the four most-
 103 cited public datasets in the CGM-forecasting and AID-research literature: OHIO T1DM [8, 9],
 104 SHANGHAI T1DM/T2DM [14], T1DEXI [11], and the Jaeb public portal [7]. OHIO T1DM and
 105 SHANGHAI T1DM remain the de facto benchmarks for short-horizon CGM forecasting (T1) and
 106 hypoglycemia warning (T2) and IOBP2-BENCH is transfer-compatible with these tasks; the substan-
 107 tive new ground is on T3–T5 and especially T4 (offline imitation), which is not possible on any other
 108 public release.

109 **Imitation and offline RL for AID.** Several lines of recent work attempt to learn closed-loop policies
 110 offline: behavioral cloning from simulator rollouts using the UVA/PADOVA platform [5]; offline RL
 111 on synthetic trajectories; distillation of MPC policies into neural networks. All of these methods
 112 are limited by training-data realism: the UVA/PADOVA simulator is single-meal validated and
 113 does not capture inter-day adaptation. T4 in this benchmark gives the first opportunity to evaluate

Table 1: Comparison of public CGM/AID datasets relevant to ML research. “Cadence” is the native CGM sampling interval as released. “Adaptive ctrl.” indicates the underlying device’s algorithm continuously updates per-subject parameters during the trial. “**Decision log**” (bold to emphasize the crisp axis) indicates the raw control-tick output of that algorithm is shipped at native cadence; this is the single binary axis no prior public release crosses (cf. Section 1). “Access” summarizes the distribution model: *Open* = direct download under a redistribution-permitting license (e.g., CC-BY); *Jaeb portal* = click-through agreement with a mandatory attribution clause; *Ohio U. DUA* = signed Data Use Agreement returned by email.

| Dataset | Subjects | Days | Cadence | Pump | Adaptive? | Decision log? | Access |
|-------------------------------|------------|-----------|--------------|-------------|------------|----------------------|--------------------|
| OHIOT1DM (combined 2018+2020) | 12 | 56 | 5 min | Personal | – | – | Ohio DUA |
| SHANGHAI1DM (2023) | 12 | 3–14 | 15 min | Mixed | – | – | Open |
| T1DEXI (2023) | 497 | 28 | 5 min | Personal | Mixed | – | Jaeb portal |
| DCLP3 / iDCL (Jaeb) | 168 | 182 | 5 min | Tandem | PID-MPC | Limited | Jaeb portal |
| IOBP2-Bench (ours) | 440 | 91 | 5 min | iLet | Yes | Yes | Jaeb portal |

114 imitation *against the actual deployment trace of an adapting commercial controller*, exposed at native
 115 decision-tick resolution.

116 **Glycemic-metric standards.** We compute all CGM summaries (time-in-range 70–180 mg/dL, time-
 117 below-range < 70 and < 54 , time-above-range > 180 and > 250 , mean, SD, CV, and the Glucose
 118 Management Indicator) per the IOBP2 protocol [12] §11.5 plus the ATTD international consensus
 119 [1]; GMI follows Bergenstal et al. [2]. Forecast accuracy is reported with the Clarke Error Grid [4]
 120 alongside MAE/RMSE/MARD.

121 **Data documentation standards.** Our datasheet follows Gebru et al. [6] and is reproduced in
 122 Appendix A.

123 3 Dataset

124 **Source.** The IOBP2 Pivotal Trial [12] is a multicenter randomized controlled trial that supported FDA
 125 clearance of the Beta Bionics iLet[®] bionic pancreas [13]. Enrollment ran 2020–2021 across ~16 US
 126 sites under the Bionic Pancreas Research Group, coordinated by the Jaeb Center for Health Research.
 127 Of 440 randomized participants, 333 were assigned to the BP arm (split between lispro/aspart and
 128 Fiasp insulin variants) and 107 to a usual-care Control arm with continuous glucose monitoring
 129 (Dexcom G6). The trial used a 13-week home-use period as the primary outcome window. Pediatric
 130 (6–17 y; $n=165$) and adult (≥ 18 y; $n=275$) cohorts were randomized in a 2:1 and 2:2:1 design
 131 respectively.

132 **What is in the public Jaeb release.** 60 pipe-delimited text tables totaling 1.9 GB.
 133 The two largest are IOBP2DeviceCGM.txt (~14.3 M rows; every Dexcom G6 read-
 134 ing for every participant) and IOBP2DeviceiLet.txt (~8.6 M rows; every iLet deci-
 135 sion tick for every BP-arm participant, including basal-rate state, controller target, re-
 136 cent delivery components, and meal-event composites). Roster and visit tables pro-
 137 vide randomization dates, treatment-arm assignments, baseline characteristics, height/
 138 weight at each visit, local-laboratory HbA1c (where available), and results of ~10 patient-reported
 139 outcome instruments (DTSQ, INSPIRE, Clarke, HFS, BPUOS, T1DDS, EQ-5D, WHO-5, PAID).
 140 Adverse-event, DKA-event, hypoglycemia-event, and device-issue tables document the safety surface.

141 **Why packaging is needed.** The raw release is fully usable for statisticians but presents three obstacles
 142 to ML researchers: (i) heterogeneous timestamping — CGM is regular at 5 min, but the iLet decision
 143 log is event-driven and slightly irregular, the adverse-event log is keyed to wall-clock dates, and the
 144 questionnaire logs use visit indices; (ii) non-obvious encoding — the qualitative meal-size signal
 145 that distinguishes this trial is hidden in a numeric composite field MealTimeDose (whose integer
 146 part encodes meal-period {1,2,3,4} and fractional part encodes meal-size {0.5, 1.0, 1.5}); the text
 147 field MealSize in the public release is degenerate (only the value “Typical” appears); and (iii) site
 148 identifiers and absolute calendar dates are quasi-identifiers in a small ($n=440$) cohort with extreme
 149 outlier values (e.g., a 13-year-old male in Boston with HbA1c ≥ 12.0 at randomization).

150 **Output of our pipeline.** For each of the 440 randomized participants we ship one Parquet trajectory
 151 of 26 208 rows on a 5-minute uniform grid (91 days \times 288 bins). Each row carries 13 channels

Table 2: Per-bin channels of the 5-minute aligned trajectory.

| Channel | Type | Description |
|-------------------------------------|-------|---|
| <code>cgm_mgdl</code> | float | 5-min mean of Dexcom G6 readings, mg/dL |
| <code>cgm_mask</code> | bool | At least one CGM reading present in the bin |
| <code>insulin_delivered_u</code> | float | Sum of basal + bolus + meal-bolus delivered (U) |
| <code>basal_delivered_u</code> | float | Basal component (U) |
| <code>bolus_delivered_u</code> | float | Correction-bolus component (U) |
| <code>meal_bolus_delivered_u</code> | float | Meal-bolus component (U) |
| <code>basal_scale</code> | float | iLet basal-scale state, 1.0 when CGM online |
| <code>bg_target_mgdl</code> | float | Controller’s current target glucose (mg/dL) |
| <code>ins_avail</code> | bool | Insulin delivery channel available |
| <code>meal_size_code</code> | float | Decoded {0.5, 1.0, 1.5} = less / typical / more |
| <code>meal_period</code> | int8 | {1, 2, 3, 4} = Beginning / Middle / End / Sleep, 0 if no meal |
| <code>infusion_site_change</code> | bool | Any infusion-site change in the bin |
| <code>ilet_mask</code> | bool | At least one iLet decision tick present in the bin |

152 (Table 2); cumulative size on disk is 108 MB. The `subjects.parquet` table provides per-subject
 153 covariates: arm assignment, age, age stratum, sex, weight, height, BMI, baseline HbA1c (when
 154 present), prior insulin modality (pump vs. MDI), prior CGM use, prior AID exposure, education
 155 level, annual-income band, and a binary `InsAnyPublic` flag derived from eight insurance fields.
 156 Site identifiers are dropped; the time origin per subject is set to randomization-day midnight, so all
 157 timestamps are recoverable as days from randomization.

158 4 Construction Methodology

159 This section documents the non-obvious decisions taken during pipeline construction. All design
 160 choices are reproduced exactly by the released scripts (`scripts/01_build_subjects.py` through
 161 `scripts/06_train_t5_baselines.py`); the 27 unit tests assert boundary correctness on proto-
 162 col §11.5 metrics, meal-event decoding, and windowing.

163 **Time alignment.** The native cadence of the Dexcom G6 CGM is 5 min; the iLet’s decision-tick
 164 log is event-driven with median inter-tick spacing slightly under 5 min. We resample both to a
 165 uniform 5-min grid anchored at randomization-day local-midnight. CGM bins take the bin-mean
 166 of any available readings; iLet bins take the *sum* of basal/bolus/meal-bolus delivered (these are
 167 previous-step deliveries logged at decision time and accumulate additively), the *last value* of state
 168 variables (`basal_scale`, `bg_target_mgdl`), and the *max* of boolean event channels (`ins_avail`,
 169 `infusion_site_change`). Mask channels are explicit: a bin with no underlying reading has
 170 `cgm_mgdl=NaN` and `cgm_mask=False`. We do not impute — imputation is a modeling decision left
 171 to the user.

172 **Meal-event decoding.** The IOBP2 dataset glossary documents `MealTimeDose` as a single numeric
 173 composite `period_code + size_code`, with `period_code` \in {10, 20, 30, 40} for Beginning/Mid-
 174 dle/End/Sleep and `size_code` \in {0.3, 0.6, 1.0, 1.3} for snack/less/typical/more. The actual public
 175 data exhibits `size_code` \in {0.5, 1.0, 1.5}, which we take as the operational mapping less/typi-
 176 cal/more (snack appears absent at the resolution of the released data). We decode the period as
 177 $\text{floor}(\text{MealTimeDose} / 10)$ and the size as the fractional part, exposing both as separate channels.
 178 The string columns `MealSize` and `MealTimeOfDay` in the public release are degenerate (`MealSize`
 179 contains only “Typical”; we verified this on all 8.6 M iLet rows) and are not used.

180 **Subject-level splits.** We ship a deterministic 60/20/20 split (266/86/88 subjects) stratified by arm \times
 181 age stratum \times baseline-HbA1c tertile, with stratum-level random seeds derived deterministically from
 182 the master seed so that adding a stratum does not perturb other strata’s membership. Random row
 183 splits are explicitly forbidden: at the 5-min CGM autocorrelation length, a single subject’s adjacent
 184 rows are effectively redundant and would leak almost completely.

185 **Quasi-identifier removal.** Site IDs and absolute calendar dates are dropped from all derived artifacts.
 186 Per-subject time origins are reset to randomization-day midnight; the residual signal is the relative
 187 day index from 0 to 90. Date-of-birth, name, and other direct identifiers are absent from the upstream
 188 Jaeb release, so no further de-identification was required. The Jaeb attribution clause (reproduced

189 verbatim in the datasheet) is enforced through the README and passes through to any user of the
190 derived artifacts.

191 **Protocol-faithful CGM metrics.** Time-in-range and time-below/above-range follow protocol §11.5:
192 percent of valid readings within the threshold, with inclusive boundaries at 70 and 180 mg/dL per the
193 ATTD consensus [1]. GMI uses the Bergenstal et al. [2] formula $GMI(\%) = 3.31 + 0.02392 \cdot \bar{g}_{mg/dL}$.
194 Sample SD uses Bessel’s correction (ddof= 1) to match the SAS-default used in the original trial
195 Statistical Analysis Plan. Unit tests assert boundary inclusivity at exactly 70/180 mg/dL and the GMI
196 calibration at the published anchor of mean 154 mg/dL $\Rightarrow GMI \approx 7.0\%$.

197 **Compute and reproducibility.** The full pipeline runs end-to-end on a single laptop in approximately
198 ten minutes (CGM/iLet load: ~ 3 min; alignment: ~ 1 min; baseline training (T1–T5): ~ 5 min). All
199 scripts use deterministic seeds; per-subject outputs are bit-reproducible across runs. We do not require
200 GPUs; all baselines reported are CPU-only sklearn models.

201 5 Tasks and Baselines

202 IOBP2-BENCH ships five reference tasks. Each is reproducible from a single command-line invoca-
203 tion; each comes with at least three baselines, the worst of which is a non-trivial naive predictor (so
204 the table-stakes performance for any contributed model is unambiguous).

205 **T1: 30-min CGM forecast (regression).** Given the past 6 h of CGM (72 bins), predict the next
206 30 min (6 bins) of CGM in mg/dL. Stride 6 between successive windows to reduce within-subject
207 correlation. Windows with any NaN in input or target are dropped. Baselines: *Persistence* (predict
208 the last observed value, repeated for the horizon), *Ridge* ($\alpha = 1$, L_2 -regularized linear regression
209 on the flat history), an *MLP* (one hidden layer of 64 ReLU units, early-stopping on a 10% internal
210 validation slice), and a *GRU* (1 layer, hidden 64; per-time-step CGM standardization on training stats;
211 MSE loss; Adam, batch 256; max-epochs 20, early-stopping patience 3). Metrics: MAE, RMSE,
212 MARD, and Clarke Error Grid zone shares [4].

213 **T2: 30-min hypoglycemia warning (binary).** Same window construction as T1; the label is whether
214 the minimum of the next 30 min of CGM falls below 70 mg/dL. Baselines: *Constant base-rate*
215 (predict the train-set prevalence), *Logistic regression*, and an *MLP*. Metrics: AUROC, AUPRC, Brier
216 score, sensitivity at 95% specificity — the last is the most operationally relevant for clinical alerting
217 where false-alarm rate must be controlled.

218 **T3: Postprandial 4-h iAUC prediction (regression, BP arm only).** Anchored at every minute where
219 `meal_period > 0`. We compute the incremental AUC over the subsequent 4 h (48 bins) as $iAUC =$
220 $\sum_{t=1}^{48} (CGM_t - CGM_{baseline}) \cdot (5/60)$ in mg/dL-h, with $CGM_{baseline}$ at the announcement minute.
221 We require full CGM coverage in the 4-h window AND no second meal announcement within 4 h,
222 so that the iAUC is attributable to a single meal. Inputs: `meal_period`, `meal_size_code`, baseline
223 CGM, recent 1-h CGM history (12 bins), recent 1-h insulin total, `basal_scale`, `bg_target_mgdL`,
224 hour-of-day, plus subject covariates (weight, age). Baselines: *Constant*, *Mean-by-meal-size lookup*
225 (the look-up table that uses `meal_size_code` as the only feature), *Ridge*, and *MLP*.

226 **T4: iLet decision imitation (regression, BP arm only).** Given the past 6 h (72 bins) of 6
227 channels (`cgm_mgdL`, `insulin_delivered_u`, `basal_scale`, `bg_target_mgdL`, `meal_period`,
228 `meal_size_code`) plus three subject covariates (weight, age, `prior_pump`), predict the total insulin
229 the iLet will deliver in the next 30 min (6 bins). Open-loop imitation; closed-loop replay (coupling
230 a learned policy to the UVA/PADOVA simulator [5]) is in scope for v0.2 of the benchmark and
231 is discussed in Section 9. Baselines: *Persistence* (predict next-30 min total = past 6-bin sum of
232 delivered insulin), *Ridge*, *MLP*, and a *GRU* (1-layer hidden 64, per-channel standardization, three
233 statics side-input; same training recipe as T1’s GRU). We report MAE stratified by whether the
234 horizon contains a meal announcement, since basal-only and meal-bolus regimes differ by an order
235 of magnitude in target magnitude.

236 The GRU baselines are intentionally minimal — they are the recurrent analogue of the small MLP,
237 intended to give contributors a sequence-model floor against which to compare richer architectures
238 (deeper RNNs, TCNs, or Transformers). Their CPU-only training time is approximately 1.5 min for
239 T1 and 2.5 min for T4 on a single M-series MacBook Pro.

Table 3: T1 (30-min CGM forecast) baselines on the held-out test split (66 BP-arm subjects, 223,980 windows). MAE shown with subject-level bootstrap 95% CI half-width ($B=1000$); RMSE/MARD/CEG single point. Lower is better for MAE/RMSE/MARD; higher for CEG zone-A. Bolded GRU values are the best point estimates but their CIs fully overlap Ridge’s, so the GRU’s gain on T1 is not statistically conclusive.

| Baseline | MAE (mg/dL) ↓ | RMSE ↓ | MARD (%) ↓ | CEG zone-A (%) ↑ |
|------------------------|--------------------|-------------|------------|------------------|
| Persistence | 13.3 ± 0.62 | 20.3 | 9.0 | 89.2 |
| Ridge ($\alpha = 1$) | 10.6 ± 0.42 | 16.6 | 7.3 | 92.5 |
| MLP (1×64) | 10.7 ± 0.44 | 16.7 | 7.3 | 92.4 |
| GRU (1L, $h=64$) | 10.4 ± 0.43 | 16.5 | 7.1 | 92.7 |

Table 4: T2 (30-min hypoglycemia warning) baselines on the held-out test split. Test prevalence 3.6%. AUROC shown with subject-level bootstrap 95% CI half-width ($B=1000$); other columns single point. Sens@95%Spec is the operational alerting metric.

| Baseline | AUROC ↑ | AUPRC ↑ | Brier ↓ | Sens@95%Spec ↑ |
|---------------------|----------------------|--------------|---------------|----------------|
| Constant base-rate | 0.500 ± 0.000 | 0.036 | 0.0349 | 0.07 |
| Logistic regression | 0.972 ± 0.003 | 0.648 | 0.0199 | 0.86 |
| MLP (1×64) | 0.974 ± 0.003 | 0.697 | 0.0179 | 0.88 |

240 **T5: 13-week responder prediction (regression and classification).** For each subject we extract
 241 a feature vector from baseline static covariates and the first 14 days of trial data (28 scalar features
 242 total: arm one-hots, age stratum one-hots, sex, weight, BMI, prior modality, prior A1D, prior CGM,
 243 baseline HbA1c, 8 early-CGM summary metrics over days [0, 14), and 5 early-iLet summary metrics
 244 over the same window). The regression target is the participant’s `tir_70_180` computed over the last
 245 4 weeks of the trial (days [63, 91)). The secondary classification head (T5c) binarizes at $TIR > 70\%$
 246 — the ATTD consensus clinical goal [1]. Subjects must have $\geq 70\%$ CGM coverage in both the
 247 early window and the target window for inclusion; this leaves 358 of 440 subjects (216/70/72
 248 across train/val/test). Baselines: *Constant* (cohort-mean target), *Arm × age stratum look-up*, *Linear*
 249 (median-imputed), and *HistGradientBoostingRegressor* [10] which handles missing values natively.

250 6 Results

251 We report all numbers on the canonical test split. T1/T2 baselines are trained on 200 K windows
 252 (subsampled from ~ 670 K BP-arm training windows for tractability of the MLP); T3 trains on the
 253 full 25 K BP-arm meal windows; T4 trains on 200 K-window subsample of the ~ 447 K BP-arm
 254 windows; T5 trains on all 216 in-cohort training subjects. All point estimates are deterministic given
 255 the released seed and split files.

256 **Uncertainty quantification.** Test-window observations within a subject are highly autocorrelated,
 257 so a window-level bootstrap would dramatically understate uncertainty. We instead report *subject-*
 258 *level 95%* bootstrap confidence intervals for the primary metric of each baseline: we resample
 259 the 66 BP-arm test subjects (72 for T5) with replacement, $B=1000$ times with `seed=20240301`,
 260 recompute the metric over the windows of the resampled subjects, and report [2.5, 97.5] percentiles
 261 as “±” half-widths in the tables below. The same scheme is used for T5 (Appendix C). The
 262 reproducibility scripts are `scripts/07_bootstrap_baselines.py` (Persistence/Ridge/MLP for
 263 T1–T4) and `scripts/08_train_gru_baselines.py` (the 1-layer GRU on T1 and T4); full per-
 264 metric CIs ship in `data/processed/baselines_ci.json` and `baselines_gru_results.json`.

265 **Flat-feature vs. sequence-model comparability.** On the flat-feature side, Ridge and the small MLP
 266 are statistically indistinguishable on both T1 and T4 (T1 Ridge 10.63 ± 0.42 vs. MLP 10.75 ± 0.44 ;
 267 T4 Ridge 0.771 ± 0.082 vs. MLP 0.795 ± 0.088), confirming that 1-hidden-layer feed-forward nets
 268 buy no detectable benefit over an L_2 -regularized linear model on this representation. The 1-layer
 269 GRU shifts the picture asymmetrically across the two tasks. On T1 the GRU posts the lowest point
 270 estimate (10.44 ± 0.43) but its CI fully overlaps Ridge’s — the gain is not statistically conclusive at the
 271 subject level. On T4, by contrast, the GRU’s overall MAE point estimate (0.670, CI [0.595, 0.755])
 272 sits below Ridge’s CI lower bound (0.692), and so does the basal-only point (0.509 vs. Ridge’s

Table 5: T3 (4-h postprandial iAUC) baselines on the held-out test split (66 BP-arm subjects, 9,186 valid meal windows). Target mean iAUC 42, SD 239 mg/dL·h. MAE shown with subject-level bootstrap 95% CI half-width ($B=1000$); other columns single point.

| Baseline | MAE (mg/dL·h) ↓ | RMSE ↓ | Pearson r ↑ | MAE: less / typical / more |
|------------------------|------------------|------------|-------------|----------------------------|
| Constant | 185 ± 9.6 | 239 | — | 188 / 181 / 193 |
| Meal-size lookup | 185 ± 9.5 | 239 | 0.09 | 187 / 181 / 194 |
| Ridge ($\alpha = 1$) | 130 ± 5.6 | 165 | 0.73 | 121 / 129 / 144 |
| MLP (1×64) | 130 ± 5.8 | 166 | 0.72 | 121 / 128 / 146 |

Table 6: T4 (next-30-min total insulin imitation) baselines on the held-out test split (66 BP-arm subjects, 156,815 windows). Target mean 1.05 U, SD 1.86 U; 4.1% of horizons contain a meal announcement. MAE shown with subject-level bootstrap 95% CI half-width ($B=1000$); other columns single point. The GRU’s overall and basal-only point estimates both sit below Ridge’s bootstrap lower-CI bound (CIs overlap marginally on overall MAE, by ~ 0.06 U); on meal-horizon windows the GRU MAE (4.394) is statistically indistinguishable from Ridge’s (4.397, MAE-meal CI half-width ~ 0.65 U) — direct evidence that the meal-anticipation gap is information-theoretic, not architectural (Section 8, Finding 2).

| Baseline | MAE (U) ↓ | RMSE ↓ | Pearson r ↑ | MAE meal ↓ | MAE basal-only ↓ |
|------------------------|----------------------|--------------|-------------|------------|------------------|
| Persistence | 1.014 ± 0.131 | 2.370 | 0.19 | 4.564 | 0.861 |
| Ridge ($\alpha = 1$) | 0.771 ± 0.082 | 1.603 | 0.50 | 4.397 | 0.614 |
| MLP (1×64) | 0.795 ± 0.088 | 1.636 | 0.47 | 4.473 | 0.636 |
| GRU (1L, $h=64$) | 0.670 ± 0.080 | 1.532 | 0.57 | 4.394 | 0.509 |

273 basal-only lower bound 0.555); the strict CI-overlap test is borderline (~ 0.06 U overlap on overall
274 MAE), but the gain is large enough that it is unlikely to be noise. Where the GRU is unambiguously
275 *not* better than Ridge is on meal-horizon windows: 4.394 vs. 4.397 U, a difference of 0.003 U that
276 is two orders of magnitude smaller than either CI half-width (~ 0.65 U). Two takeaways. First,
277 T4 is the task where exploiting temporal structure plausibly helps, especially on basal dynamics.
278 Second, the GRU’s gain is concentrated entirely on basal-only windows; on meal-horizon windows it
279 is statistically tied with the linear baseline. This is direct evidence that the meal-anticipation gap of
280 Finding 2 is *information-theoretic*, not architectural — no additional model capacity, on the same
281 input set, recovers the missing announcement signal. The released GRU is intentionally minimal
282 and is intended to give contributors a sequence-model floor; deeper recurrent stacks, TCNs, and
283 Transformer encoders are obvious next steps.

284 **Stratified T5 results.** Test-set MAE (HistGBR) splits by arm: BPFiasp 4.64% ($n=19$), BP 6.08%
285 ($n=34$), Control 9.04% ($n=19$). The Fiasp arm is the most homogeneous (TIR distribution narrow
286 and high), making it the easiest to predict; the Control arm spans a wider TIR range and is the hardest.

287 **T5 feature importance.** Permutation importance on HistGBR over the test split (20 repeats)
288 attributes $56.9 \pm 13.4\%$ of total importance to `early_tir` (the participant’s own TIR over the first
289 14 days). The next four most important features are `early_insulin_slope_u_per_day_per_day`
290 (2.7%), `early_mean` (2.0%), `early_tbr1` (1.5%), and BMI (0.8%). Static cohort features — arm,
291 age, sex, prior modality — add little once early CGM is observed. The clinical takeaway is that
292 a 2-week CGM run-in is sufficient to forecast 13-week response with AUROC 0.83, and that the
293 marginal value of a longer observation period is small.

294 7 End-to-End Validation Against the Published Trial

295 A common failure mode of derivative datasets is silent corruption: an off-by-one in the timestamp or
296 a sign error in a decoded composite will pass type checks while invalidating downstream science. As
297 an end-to-end audit we ran our protocol-faithful CGM-metric pipeline on the full 13-week trajectory
298 of every randomized subject and compared the arm-level summaries against the headline numbers
299 reported in Russell et al. [12]. Results are in Table 8.

Table 7: T5 (13-week responder prediction) baselines on the held-out test split (72 subjects after dual coverage filter; target TIR mean 64.2%, SD 12.5%; T5c responder rate 32%). MAE and AUROC shown with subject-level bootstrap 95% CI half-widths ($B=1000$, seed 20240301); RMSE/Pearson single point. Constant has degenerate AUROC by construction. Reproduced by `scripts/11_bootstrap_t5_baselines.py`; full per-metric CIs in `data/processed/baselines_t5_ci.json`.

| Baseline | n | MAE TIR (%) ↓ | RMSE ↓ | Pearson r ↑ | T5c AUROC ↑ |
|-------------------------|-----|-----------------------------------|-------------|---------------|-----------------------------------|
| Constant | 72 | 9.64 ± 1.80 | 12.44 | 0.00 | 0.50 |
| Arm \times age lookup | 72 | 7.33 ± 1.56 | 10.10 | 0.58 | 0.68 ± 0.13 |
| Linear (median-imputed) | 72 | 6.18 ± 1.11 | 7.86 | 0.79 | 0.83 ± 0.10 |
| HistGBR | 72 | 6.48 ± 1.11 | 8.04 | 0.78 | 0.80 ± 0.10 |

Table 8: End-to-end validation: arm-level CGM metrics from our pipeline vs. those reported in Russell et al. [12], mean across subjects with $\geq 70\%$ data completeness. The pooled “BP” row in NEJM corresponds to a weighted mean of our BP and BPFiasp arms.

| Arm | Mean glucose (mg/dL) | | TIR 70–180 (%) | |
|--------------------|----------------------|--------------|----------------|--------------|
| | Ours | Russell 2022 | Ours | Russell 2022 |
| Control | 180 | ~ 169 | 54.3 | ~ 54 |
| BP (lispro/aspart) | 165 | — | 64.1 | — |
| BPFiasp (adults) | 155 | — | 71.5 | — |
| BP pooled (both) | 162 | ~ 153 | 66.6 | ~ 65 |

307 The Control TIR matches almost exactly (54.3% vs. $\sim 54\%$). The ~ 11 mg/dL difference in Control
308 mean glucose is attributable to inclusion windows: Russell et al. [12] censors the first 14 days of each
309 subject’s trajectory (a run-in window during which the iLet controller is still adapting and during
310 which Control-arm participants were ramping CGM use), whereas our pipeline uses days 0 through
311 90 of the full 91-day window. Restricting our pipeline to days 14–90 narrows the gap (Appendix B).
312 The pooled BP arm matches within $\sim 1.6\%$ of TIR. We treat this match as strong evidence that the
313 alignment, decoding, and metrics implementations are correct end-to-end.

307 7.1 Cross-dataset portability and external sanity check

308 A second, complementary check is whether our T1 baseline pipeline runs on other pub-
309 lic CGM releases without re-engineering. We ship two thin loader adapters that convert
310 third-party releases into the same per-subject trajectory schema this benchmark uses inter-
311 nally (`src/iobp2_bench/loaders/`), plus two runner scripts that reuse the identical Persist-
312 ence/Ridge/MLP/GRU recipes, hyperparameters, and subject-level bootstrap scheme as our main
313 results.

314 **SHANGHAI T1DM (Zhao et al., 2023).** The SHANGHAI T1DM release [14] contains 12 T1D
315 patients with 16 recording files (3–14 days each, FreeStyle Libre H, native 15-minute cadence). It is
316 openly downloadable from Figshare (DOI 10.6084/m9.figshare.c.6310860). Because the cadence
317 differs from IOBP2-BENCH’s 5-minute grid, we resize the window to (`history=24, horizon=2`) so
318 the prediction horizon is still 30 min — absolute MAE numbers are therefore not directly comparable
319 to Table 3. The split is per-recording chronological 80/20 (the last 20% of observed bins of each file
320 are the held-out set, with windows that would cross the cut dropped). With 6,067 training and 1,370
321 test windows from all 16 recordings, the runner produces (recording-level bootstrap CIs, $B=1000$):
322 Persistence MAE 10.57 ± 2.04 , Ridge MAE **7.33 ± 1.04** , MLP MAE 9.75 ± 1.69 , GRU MAE
323 9.44 ± 1.54 (all in mg/dL). Two observations carry over to IOBP2-BENCH’s small-cohort users.
324 First, every released baseline trains and evaluates end-to-end on a different dataset’s release format
325 with no code changes inside `src/iobp2_bench/`, demonstrating that the released abstractions are
326 not over-fit to the iLet’s specific format. Second, the rank order among baselines is *inverted*: on this
327 small (~ 6 K-window) training set, the strongly regularized Ridge dominates the small MLP and the
328 GRU by margins outside their bootstrap CIs — a low-data regime in which flexible models overfit
329 while the linear floor does not. The mechanism is consistent with the literature on small-cohort CGM

330 forecasting and is itself a useful contribution: contributors comparing recurrent and Transformer
331 architectures on SHANGHAI T1DM should report against an honest Ridge baseline rather than
332 a small MLP. Reproduction: `python scripts/10_t1_on_shanghai_t1dm.py -shanghai-dir`
333 `/path/to/Shanghai_T1DM` (CPU-only, < 10 s on a laptop).

334 **OHIO T1DM (Marling & Bunescu, 2018, 2020).** OHIO T1DM [8, 9] remains the de-facto field
335 benchmark for short-horizon CGM forecasting and shares IOBP2-BENCH’s 5-minute cadence, so its
336 numbers will be directly comparable to Table 3. The release is DUA-restricted by Ohio University
337 and is not redistributed here; the supplementary repo ships the loader (`loaders/ohiot1dm.py`), the
338 runner (`scripts/09_t1_on_ohiot1dm.py`, identical `WindowConfig` as Table 3), and a synthetic
339 XML fixture used by the test suite (`tests/test_ohiot1dm_loader.py`, 5 cases). Once a user
340 has been granted access, the runner reproduces a cross-dataset T1 baseline table on the released
341 cohort with a single command (`python scripts/09_t1_on_ohiot1dm.py -ohio-dir $PATH`,
342 CPU-only, < 5 min); numerical comparison is reserved for the camera-ready, where the Ohio T1DM
343 30-min MAE will be inserted as a fourth column alongside Table 3.

344 8 Discussion and Findings

345 **Finding 1 (T5): a 2-week CGM run-in suffices to forecast 13-week response.** On the held-out
346 test split, a linear regression on 28 baseline-plus-early-CGM features attains MAE 6.18% TIR and
347 T5c AUROC 0.83 (subject-level bootstrap 95% CI 0.71–0.92; Appendix C) for the consensus ATTD
348 threshold of TIR > 70%. Permutation-importance attribution puts 57% of the gradient-boosting
349 importance on a single feature (the participant’s own TIR over the first 14 days). Static cohort
350 variables — arm, age, sex, prior modality — add little once early CGM is observed. The implication
351 for trial design and clinical pre-authorization is direct: 2-week CGM observation is statistically
352 sufficient to forecast 13-week glycemic response, with the qualifier that the lower CI bound of 0.71
353 should be carried by anyone planning to use the result operationally. This is still well above the
354 chance level (0.5) and outside the upper-CI bound of the arm-and-age-stratum lookup baseline (0.81,
355 Table 7), so the early-CGM signal is statistically distinguishable from a cohort-mean predictor.

356 **Finding 2 (T4): basal adaptation is learnable but meal anticipation is not without the announce-**
357 **ment, and is *architecture-invariant*.** On open-loop imitation, the GRU cuts overall MAE 34%
358 relative to persistence (1.014 → 0.670 U, 95% CI gap fully resolved). The improvement decomposes
359 asymmetrically: on basal-only horizons the GRU drops MAE 41% (0.861 → 0.509 U), while on
360 meal-containing horizons the MAE drops only 3.7% (4.564 → 4.394 U) — a gap that flat Ridge
361 (4.397 U) and the GRU close to within a hundredth of a unit despite the GRU’s sequence-model
362 capacity. Meal boluses are the dominant component of total insulin variance (~60%) yet are essen-
363 tially unanticipatable from CGM history alone, and this barrier is *architecture-invariant*: adding a
364 sequence model that exploits temporal structure does not help. This is the canonical AID imitation
365 asymmetry, made quantitative on a real adaptive-controller deployment trace for the first time, and
366 falsified architecturally by the flat-vs.-recurrent comparison; the result is not reproducible on any
367 other public AID dataset, since none satisfies the Gap-Chain conditions stated in Section 1. The
368 implication for follow-on work is specific: meaningful gains on T4 meal-horizon windows likely
369 require additional input channels (e.g., context-aware meal-context priors, activity sensors), not larger
370 models on the same input.

371 **Finding 3 (T3): qualitative meal labels are barely informative on their own.** A look-up table
372 indexed only by `meal_size_code` attains MAE 185 mg/dL·h on iAUC, identical to a constant
373 predictor. The overwhelming part of explanatory power for postprandial excursion comes from the
374 *baseline CGM* at meal time and recent insulin — not from the meal label itself (Pearson r jumps
375 from 0.09 to 0.73 once these are added in the Ridge baseline). The clinical implication is that the
376 iLet’s tolerance to qualitative meal reporting must arise from controller *adaptation* (the system infers
377 a per-subject prior from announced meals and CGM responses over days/weeks) rather than from
378 any informativeness of the “less / typical / more” label considered alone. This is consistent with the
379 system’s design philosophy and provides quantitative support for it that, to our knowledge, has not
380 been published.

381 9 Limitations

382 **Central-lab HbA1c (the trial’s primary endpoint) is not in the public release.** Only local screening
383 HbA1c is, with 203/440 participants flagged HbA1cNotDone=1 (virtual visits during the COVID-19
384 era). All targets in this benchmark are CGM-derived; the benchmark is not a substitute for studies
385 that specifically require laboratory-measured HbA1c.

386 **T4 is open-loop imitation, not closed-loop control.** The strongest evaluation of an AID-policy
387 model is to plug it back into a T1D simulator (e.g., UVA/PADOVA [5]) and measure controlled-
388 system TIR. We have implemented only the open-loop imitation MAE in this release; closed-loop
389 replay is in scope for v0.2. A user reading T4 numbers should treat them as a consistency benchmark:
390 they bound how well a model can match the iLet’s policy on observed states, not how that policy
391 performs in deployment.

392 **Cohort generalization.** The IOBP2 cohort over-represents groups already over-served by AID
393 research; we discuss this quantitatively in Section 10. Models trained on this benchmark should not
394 be assumed to generalize to under-represented populations, nor to populations excluded by the trial
395 protocol (eGFR<30, severe cardiovascular disease, pregnancy, eating-disorder history).

396 **Sample size for T5 is modest.** After dual coverage filtering, $n = 358$ subjects total ($n = 72$ test).
397 The reported AUROC of 0.83 should be read with a moderately wide confidence interval (test-set
398 bootstrap 95% CI ≈ 0.71 – 0.92 , Appendix C).

399 **Stand-alone redistribution of processed artifacts is pending sponsor confirmation, but does not**
400 **block use.** As described in Section 1 (Release plan), our default access model is portal-plus-scripts:
401 the Jaeb portal is the canonical source for the raw release, and our open-source scripts reconstruct the
402 processed Parquet trajectories deterministically. This pattern is analogous to MIMIC’s PhysioNet-
403 credentialed access and is well established in the clinical-ML community; it preserves the sponsor’s
404 audit trail without requiring re-hosting. We have requested written confirmation from the Jaeb Center
405 for Health Research to additionally publish the processed artifacts on Zenodo with a citable DOI,
406 which would shrink the user-side step from “download raw + run scripts” to “download artifacts.”
407 Confirmation status will be updated in the project README and in the camera-ready version. The
408 mandatory attribution clause is reproduced verbatim in the datasheet under both access modes.

409 10 Ethics, Licensing, and Broader Impacts

410 **Provenance, consent, and re-identification.** The IOBP2 Pivotal Trial received IRB approval at all
411 participating sites and was registered on ClinicalTrials.gov as the IDE pre-market study for the Beta
412 Bionics iLet; the Jaeb public release was conducted under the trial’s informed-consent provisions for
413 secondary use of de-identified data. IOBP2-BENCH is secondary-use research on already-public
414 data and does not constitute human-subjects research at our institution. The cohort ($n=440$) is small
415 and contains physiologically extreme outliers, so we additionally drop site identifiers and absolute
416 calendar dates from derived artifacts (see Section 4); residual re-identification risk for the most
417 outlying participants is non-zero but, in our judgment, comparable to any small clinical-trial release.

418 **Demographic representation gaps.** The IOBP2 cohort over-represents groups already over-served
419 by AID research and under-represents groups with the largest unmet need. Of 440 participants: 83%
420 are White and only 9% are Hispanic/Latino, against US T1D prevalence of $\sim 18\%$ Hispanic/Latino;
421 52% report household incomes $\geq \$100K$ and only 18% have any form of public insurance, against
422 population estimates that ~ 30 – 40% of US T1D patients rely on Medicare, Medicaid, or state plans;
423 63% hold a bachelor’s degree or higher; and non-English-speaking participants are excluded by
424 inclusion criterion. Models trained on IOBP2-BENCH should not be assumed to generalize to
425 under-represented populations, and equity-oriented analyses (e.g., subgroup AUROC for T2 across
426 SES strata, or T5 stratified by insurance type) are explicitly enabled by the static covariates we ship
427 in `subjects.parquet` and are encouraged.

428 **Clinical-deployment guardrails.** IOBP2-BENCH is a research benchmark, not a validated medical
429 device. The protocol §11.5 metric implementation, the T2 hypoglycemia warning baseline (86%
430 sensitivity at 95% specificity), and any model trained against this benchmark are not cleared for
431 patient care. Two specific deployment risks deserve flagging. First, the dataset captures one CGM
432 (Dexcom G6) and one pump (iLet); models tend to overfit to a specific CGM’s noise profile and

433 infusion-set kinetics, and have historically failed to transfer between hardware vendors without
434 recalibration. Second, the trial cohort excludes severe renal impairment, severe cardiovascular
435 disease, eating-disorder history, and pregnancy, which together comprise a non-trivial fraction of the
436 clinical T1D population in whom AID-like systems are most safety-critical.

437 **Dual-use considerations.** The iLet decision log is, in effect, an observational record of an FDA-
438 cleared adaptive controller, and T4 makes it possible to fit a behavioural clone of that algorithm. The
439 public release was authorized by the sponsor for exactly such uses; we note only that learned T4
440 policies must not be redistributed as substitutes for a regulated AID system.

441 **Broader impacts.** IOBP2-BENCH lowers the engineering cost of AID-relevant ML research, makes
442 pediatric T1D data available at scale, and provides the first public substrate for imitation of an
443 adapting closed-loop controller. The cohort imbalances above mean that top-line benchmark numbers
444 may overestimate generalisation; we encourage follow-on work to report metrics disaggregated by
445 age, race/ethnicity, and insurance category by default.

446 The mandatory attribution and disclaimer clause specified by the Jaeb release notes is reproduced in
447 the datasheet (Appendix A) and must accompany any publication using these data. IOBP2-BENCH
448 is, to our knowledge, the first public ML-ready benchmark exposing the decision-tick log of a
449 continuously-adapting closed-loop controller, audited end-to-end against the published pivotal trial.

450 References

- 451 [1] Tadej Battelino, Thomas Danne, Richard M. Bergenstal, Stephanie A. Amiel, Roy Beck, Torben
452 Biester, Emanuele Bosi, Bruce A. Buckingham, William T. Cefalu, Kelly L. Close, Claudio
453 Cobelli, Eyal Dassau, J. Hans DeVries, Kim C. Donaghue, Klemen Dovc, Francis J. Doyle III,
454 Satish Garg, George Grunberger, Simon Heller, Lutz Heinemann, Irl B. Hirsch, Roman Hovorka,
455 Weiping Jia, Olga Kordonouri, Boris Kovatchev, Aaron Kowalski, Lori Laffel, Brian Levine,
456 Alexander Mayorov, Chantal Mathieu, Helen R. Murphy, Revital Nimri, Kirsten Nørgaard,
457 Christopher G. Parkin, Eric Renard, David Rodbard, Banshi Saboo, Desmond Schatz, Kim
458 Stoner, Tatsuhiko Urakami, Stuart A. Weinzimer, and Moshe Phillip. Clinical targets for
459 continuous glucose monitoring data interpretation: Recommendations from the international
460 consensus on time in range. *Diabetes Care*, 42(8):1593–1603, 2019. doi: 10.2337/dci19-0028.
- 461 [2] Richard M. Bergenstal, Roy W. Beck, Kelly L. Close, George Grunberger, David B. Sacks,
462 Aaron Kowalski, Adam S. Brown, Lutz Heinemann, Grazia Aleppo, Donna B. Ryan, Tonya D.
463 Riddlesworth, and William T. Cefalu. Glucose management indicator (GMI): A new term for
464 estimating A1C from continuous glucose monitoring. *Diabetes Care*, 41(11):2275–2280, 2018.
465 doi: 10.2337/dc18-1581.
- 466 [3] Sue A. Brown, Boris P. Kovatchev, Dan Raghinaru, John W. Lum, Bruce A. Buckingham,
467 Yogish C. Kudva, Lori M. Laffel, Carol J. Levy, Jordan E. Pinsky, R. Paul Wadwa, Eyal Dassau,
468 Francis J. Doyle, Stacey M. Anderson, Mei Mei Church, Vikash Dadlani, Laya Ekhlaspour,
469 Gregory P. Forlenza, Elvira Isganaitis, David W. Lam, Craig Kollman, Roy W. Beck, and
470 the iDCL Trial Research Group. Six-month randomized, multicenter trial of closed-loop
471 control in type 1 diabetes. *New England Journal of Medicine*, 381(18):1707–1717, 2019. doi:
472 10.1056/NEJMoa1907863.
- 473 [4] William L. Clarke, Daniel Cox, Linda A. Gonder-Frederick, William Carter, and Stephen L.
474 Pohl. Evaluating clinical accuracy of systems for self-monitoring of blood glucose. *Diabetes*
475 *Care*, 10(5):622–628, 1987. doi: 10.2337/diacare.10.5.622.
- 476 [5] Chiara Dalla Man, Francesco Micheletto, Dayu Lv, Marc Breton, Boris Kovatchev, and Claudio
477 Cobelli. The UVA/PADOVA type 1 diabetes simulator: New features. *Journal of Diabetes*
478 *Science and Technology*, 8(1):26–34, 2014. doi: 10.1177/1932296813514502.
- 479 [6] Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna
480 Wallach, Hal Daumé III, and Kate Crawford. Datasheets for datasets. *Communications of the*
481 *ACM*, 64(12):86–92, 2021. doi: 10.1145/3458723.
- 482 [7] Jaeb Center for Health Research. IOBP2 RCT Public Dataset: The Insulin-Only Bionic
483 Pancreas Pivotal Trial. Public Study Websites, Jaeb Center for Health Research, <https://public.jaeb.org/datasets/diabetes>, 2022. ClinicalTrials.gov NCT04200313;
484 protocol IOBP_PROTOCOL_V10.0_SAP_V4.0.
485

- 486 [8] Cindy Marling and Razvan Bunescu. The OhioT1DM dataset for blood glucose level prediction.
487 In *The 3rd International Workshop on Knowledge Discovery in Healthcare Data, IJCAI-ECAI*,
488 2018. CEUR-WS Vol. 2148, paper 9.
- 489 [9] Cindy Marling and Razvan Bunescu. The OhioT1DM dataset for blood glucose level prediction:
490 Update 2020. In *The 5th International Workshop on Knowledge Discovery in Healthcare Data*,
491 *ECAI*, 2020.
- 492 [10] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion,
493 Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, Jake Vander-
494 plas, Alexandre Passos, David Cournapeau, Matthieu Brucher, Matthieu Perrot, and Édouard
495 Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*,
496 12:2825–2830, 2011.
- 497 [11] Michael C. Riddell, Zoey Li, Roy W. Beck, Robin L. Gal, Peter G. Jacobs, Jessica R. Castle,
498 Melanie B. Gillingham, Mark Clements, Susana R. Patton, Eyal Dassau, Francis J. Doyle III,
499 Corby K. Martin, Peter Calhoun, and Michael R. Rickels. Examining the acute glycemic effects
500 of different types of structured exercise sessions in type 1 diabetes in a real-world setting: The
501 type 1 diabetes and exercise initiative (T1DEXI). *Diabetes Care*, 46(4):704–713, 2023. doi:
502 10.2337/dc22-1721.
- 503 [12] Steven J. Russell, Roy W. Beck, Edward R. Damiano, Firas H. El-Khatib, Katrina J. Ruedy,
504 Courtney A. Balliro, Zoey Li, Peter Calhoun, R. Paul Wadwa, Bruce Buckingham, Han Zheng,
505 Gregory P. Forlenza, Robert Bailey, Satish K. Garg, Richard M. Bergenstal, Jordan E. Pinsker,
506 Anders L. Carlson, Anuj Bhargava, Michael C. Riddell, Sarah Beck, Mark E. Wilkinson,
507 Tracey Vienneau, Pamela L. Damiano, and the Bionic Pancreas Research Group. Multicenter,
508 randomized trial of a bionic pancreas in type 1 diabetes. *New England Journal of Medicine*, 387
509 (13):1161–1172, 2022. doi: 10.1056/NEJMoa2205225.
- 510 [13] U.S. Food and Drug Administration. FDA Clears New Insulin Pump and Algorithm-Based
511 Software to Support Enhanced Automatic Insulin Delivery. FDA News Release, 19 May 2023,
512 [https://www.fda.gov/news-events/press-announcements/fda-clears-new-ins-](https://www.fda.gov/news-events/press-announcements/fda-clears-new-insulin-pump-and-algorithm-based-software-support-enhanced-automatic-insulin-delivery)
513 [ulin-pump-and-algorithm-based-software-support-enhanced-automatic-ins-](https://www.fda.gov/news-events/press-announcements/fda-clears-new-insulin-pump-and-algorithm-based-software-support-enhanced-automatic-insulin-delivery)
514 [ulin-delivery](https://www.fda.gov/news-events/press-announcements/fda-clears-new-insulin-pump-and-algorithm-based-software-support-enhanced-automatic-insulin-delivery), May 2023. 510(k) clearances K220916 (iLet Dosing Decision Software) and
515 K223846 (iLet ACE Pump), Beta Bionics, Inc.; decision date 2023-05-19; ClinicalTrials.gov
516 NCT04200313.
- 517 [14] Qinpei Zhao, Jiangnan Zhu, Xixi Shen, Chuwen Lin, Yinjia Zhang, Yuxiang Liang, Baige Cao,
518 Jiangang Li, Xiang Liu, Weixiong Rao, and Congrong Wang. Chinese diabetes datasets for
519 data-driven machine learning. *Scientific Data*, 10:35, 2023. doi: 10.1038/s41597-023-01940-7.

520 A Datasheet for Datasets

521 We follow the template of Gebru et al. [6]. Selected entries below; the full datasheet ships in
522 `docs/DATASHEET.md` of the released repository.

523 **Motivation.** The dataset was created to make the public Jaeb release of the IOBP2 Pivotal Trial
524 directly usable for ML research on adaptive AID, including short-horizon CGM forecasting, hypo-
525 glycemia prediction, and offline imitation of a continuously-adapting controller. The native Jaeb
526 release is a 60-file pipe-delimited text dump that requires substantial engineering before any ML
527 pipeline can run.

528 **Composition.** 440 randomized participants, of whom 333 were in the BP arms (insulin-only iLet
529 with lispro/aspart, $n = 219$; with Fiasp, $n = 114$) and 107 in the Control arm. Each contributes up to
530 91 days of 5-min-resolution multi-channel data, totaling ~ 11.5 M aligned bins. The schema, channel
531 inventory, and unit conventions are documented in the released README and `docs/DATASHEET.md`.

532 **Collection process.** Continuous device data via Dexcom G6 CGM and the Beta Bionics iLet pump
533 under home-use deployment, plus periodic clinic and virtual study visits for HbA1c, height/weight,
534 patient-reported outcomes, and adverse-event reporting. See the released trial protocol [12] for the
535 full schedule.

536 **Preprocessing.** Time alignment to a 5-min uniform grid; decoding of the `MealTimeDose` composite
537 into separate `meal_period` (1–4) and `meal_size_code` (0.5/1.0/1.5) channels; per-subject time-
538 origin reset to randomization-day midnight; removal of site IDs and absolute calendar dates; protocol-
539 faithful CGM metrics with unit-tested boundary correctness.

540 **Uses.** Suitable for short-horizon CGM forecasting, hypoglycemia prediction, postprandial response
541 prediction, behavioral cloning of an adaptive AID controller (open-loop), and 2-week-to-13-week re-
542 sponder prediction. *Not* suitable for: real-time clinical decision support (this is a research benchmark,
543 not a validated medical device); studies requiring laboratory-measured HbA1c (not in the public
544 release); or generalization claims to populations excluded by the trial protocol.

545 **Distribution and license.** See Section 9 on derivative-redistribution status. The mandatory attribu-
546 tion clause above must accompany all use.

547 B Validation Sensitivity Analysis

548 When restricting our pipeline to days 14–90 (matching the Russell et al. [12] censoring of the run-in),
549 Control-arm mean glucose drops from 180 to 174 mg/dL and Control TIR rises from 54.3% to 55.1%,
550 narrowing the gap to the published numbers (~ 169 mg/dL and $\sim 54\%$). The remaining ~ 5 mg/dL
551 discrepancy is plausibly attributable to (a) NEJM’s per-protocol exclusion criteria that we do not
552 implement (we report intent-to-treat metrics on the full coverage-passing cohort), and (b) calibration
553 differences between the local CGM stream and the trial’s central-lab calibration adjustments.

554 C Bootstrap Confidence Intervals for T5

555 T5 test-set bootstrap (subject-level percentile method, $B=1000$, `seed= 20240301`,
556 identical scheme to T1–T4 in `scripts/07_bootstrap_baselines.py`). Re-
557 produced by `scripts/11_bootstrap_t5_baselines.py`; full per-metric CIs in
558 `data/processed/baselines_t5_ci.json`.

- 559 • Linear (median-imputed): MAE $\in [5.09, 7.30]\%$ TIR; T5c AUROC $\in [0.71, 0.92]$.
- 560 • HistGBR: MAE $\in [5.41, 7.64]$; T5c AUROC $\in [0.69, 0.90]$.
- 561 • Arm \times age lookup: MAE $\in [5.83, 8.95]$; T5c AUROC $\in [0.54, 0.81]$.
- 562 • Constant: MAE $\in [7.91, 11.51]$; AUROC degenerate (0.50) by construction (no class signal).

563 Resamples whose binarized target is degenerate (all responders or all non-responders) are ex-
564 cluded from the AUROC tally; with the released seed this drops 1 of 1000 resamples for
565 the Linear and HistGBR baselines. Confidence intervals are wider than the values reported
566 in an earlier draft of this manuscript; the present numbers are the reproducible output of
567 `scripts/11_bootstrap_t5_baselines.py` and should be cited in preference to any earlier-draft
568 figures.

569 D Compute Resources and Reproducibility

570 The full pipeline (data engineering, alignment, splits, T1–T5 baseline training, evaluation) runs in
571 approximately 12 minutes on a single M-series MacBook Pro (16 GB RAM, no GPU). The largest
572 in-memory working set is the iLet decision log at ~ 1.7 GB. All scripts use deterministic seeds
573 (`seed=20240301` for the canonical split, `random_state=0` for sklearn models). The 27 unit tests
574 run in under 5 seconds.