# Predicting T Cell Quality During Manufacturing Through an Artificial Intelligence-based Integrative Multi-Omics Analytical Platform

Valerie Y. Odeh-Couvertier, Nathan J. Dwarshuis, Maxwell B. Colonna, Bruce L. Levine, Arthur S. Edison, Theresa Kotanchek, Krishnendu Roy, and Wandaliz Torres-Garcia*

**Author information**

Valerie Y. Odeh-Couvertier, Nathan J. Dwarshuis, and Maxwell B. Colonna: These authors contributed equally to this work.

Affiliations

**Department of Industrial Engineering, University of Puerto Rico Mayagüez, Mayagüez, PR, 00681, USA**

24      Valerie Y. Odeh-Couvertier & Wandaliz Torres-Garcia

25      **The Wallace H. Coulter Department of Biomedical Engineering, Georgia Institute of**

26      **Technology, Atlanta, GA, 30318, USA**

27      Nathan J. Dwarshuis & Krishnendu Roy

28      **Departments of Genetics and Biochemistry & Molecular Biology, Complex Carbohydrate**

29      **Research Center, University of Georgia, Athens, GA, 30602, USA**

30      Maxwell B. Colonna & Arthur S. Edison

31      **Center for Cellular Immunotherapies, Perelman School of Medicine, University of**

32      **Pennsylvania, Philadelphia, PA, 19104, USA**

33      Bruce L. Levine

34      **Evolved Analytics LLC, Rancho Santa Fe, CA, USA**

35      Theresa Kotanchek

36

37      **Contributions**

38      N.J.D. designed and performed the T cell culturing experiments and measured cytokine profiles.

39      M.B.C. and A.S.E. designed and performed the NMR media analysis. V.Y.O. and W.T. executed

40      the machine learning techniques and integrated computational tools in the workflow. T.K.

41      implemented optimization and predictive analysis using DataModeler. V.Y.O., N.J.D., M.B.C.,

42      B.L.L., A.S.E., T.K., K.R., and W.T. interpreted the data and results. All authors contributed to

43      the writing, revising, and editing of the manuscript.

44

45      **Corresponding author**

46      Correspondence to Wandaliz Torres-García*.

47

48 **Ethics declarations**

49 Competing Interests

50 B.L.L. declares financial interest intellectual property and patents in the field of cell and gene

51 therapy (University of Pennsylvania Alliance with Novartis, licensing, and royalty fees). B.L.L. is

52 a consultant for Novartis, Terumo, and Lilly Asia Ventures and he is part of the Scientific Advisory

53 Board for Avectas, Brammer Bio/TF Viral Vector Services, Immuneel, Incysus, Ori Biotech, and

54 Vycellix. Moreover, B.L.L. is the co-founder and equity holder Tmunity Therapeutics and all of

55 his conflict of interest is managed in accordance with University of Pennsylvania policy and

56 oversight. T.K. is the Chief Executive Officer of Evolved Analytics, LLC. The remaining authors

57 declare no competing interests. K.R declares consulting, intellectual property, and patents in cell

58 and gene therapy. K.R. is a consultant to Terumo, Merck. LEK consulting, Mubadala Ventures,

59 Anzu Partners, Decibio, and Clearview Healthcare Partners. K.R. serves on the advisory board of

60 the MIT-Singapore Cell therapy Partnership.

61 **Abstract**

62 Large-scale, reproducible manufacturing of therapeutic cells with consistently high quality is vital

63 for translation to clinically effective and widely accessible cell therapies. However, the biological

64 and logistical complexity of manufacturing a living product, including challenges associated with

65 their inherent variability and uncertainties of process parameters, currently make it difficult to

66 achieve predictable cell-product quality. Using a degradable microscaffold-based T cell process as

67 an example, we developed an Artificial Intelligence (AI)-driven experimental-computational

68 platform to identify a set of critical process parameters (CPP) and critical quality attributes (CQA)

69 from heterogeneous, high dimensional, time-dependent multi-omics data, measurable during early

70 stages of manufacturing and predictive of end-of-manufacturing product quality. Sequential,

71 Design-of-Experiment (DOE)-based studies, coupled with an agnostic machine-learning

72 framework, were used to extract feature combinations from media assessment that were highly

73 predictive of total live CD4$^+$ and CD8$^+$ naïve and central memory (CD63L$^+$CCR7$^+$) T cells and

74 their ratio in the end-product. This computational workflow could be broadly applied to any cell

75 therapy and provide a roadmap for discovering CQAs and CPPs in cell manufacturing.

76 **Introduction**

77 T cell-based immunotherapies have received great interest from clinicians and industry due to their

78 potential to treat, and often functionally cure some cancers and their potential applicability in many

79 other diseases[1,2]. Since 2017, four genetically modified autologous Chimeric Antigen Receptor

80 (CAR) T cell therapies (*Yescarta*[TM], *Kymriah*[TM], *Tecartus*[TM], *Breyanzi*[®]) have received FDA

81 approval to treat certain B-cell malignancies. Despite these successes, CAR-T cell therapies are

82 constrained by poorly-understood manufacturing processes that are time-intensive, expensive, and

83 difficult to scale[3,4] with a lack of methods and tools to predict product quality during manufacturing

84 and identify product Critical Quality Attributes (CQAs) and the associated Critical Process

85 Parameters (CPPs).

86 Translating laboratory-scale T cell expansion experiments into a large-scale manufacturing

87 process is hindered by the incomplete understanding of cell properties and how they are affected

88 by process variables, lack of detailed characterization, and high variability of materials during

89 manufacturing[5]. These challenges of manufacturing a "living product" are further magnified since

90 current chemistry, manufacturing, and control (CMC), analytics, regulations, and product-

91 specifications are designed for conventional chemical and biopharmaceutical manufacturing

92 systems[6]. This underscores the need to develop innovative tools, methods, and standards to ensure

93 appropriate quality controls, and new strategies involving quality by design (QbD) and good

94 manufacturing practices (GMP) for cell-based therapies[7–9]. The intricate manufacturing process

95 for T cells and other cell therapies must be deeply assessed and appropriately controlled to ensure

96 scalability, predictability, and a high-quality manufacturing process at the most reasonable cost. A

97 key step for reaching this goal is to identify putative CQAs and CPPs early in the manufacturing

98 process that can predict the quality of the manufactured cell-therapy product. We hypothesized

99    that rigorous characterization of process parameters along with longitudinal measurements of cell-

100    secreted cytokine, chemokine, and metabolites from the culture media early during manufacturing

101    will allow us to develop an AI-based mathematical-computational framework for the identification

102    of multivariate parameters that are predictive of the end-of-manufacturing product phenotypes.

103    Characterization studies of approved autologous anti-CD19 CAR-T cell therapies have recently

104    revealed initial sets of candidate quality attributes, i.e. percent transduction, vector copy number,

105    and interferon-γ production for Axicabtagene ciloleucel (Yescarta[TM])[10] while CAR expression and

106    release of interferon-γ are a few of those identified for Tisagenlecleucel (Kymriah[TM])[11]. Many of

107    these attributes are calculated as endpoint responses and thus a deeper understanding of the cell

108    growth process impacted by starting conditions and performance during their manufacturing is

109    essential. Hence, CQAs that enable early monitoring through real-time process measurements such

110    as multi-omics cell characterization can overcome current challenges in assessing product

111    consistency. Yet, the computational complexity of dealing with the heterogeneity and multivariate

112    nature of multi-omics measurements to characterize T cell quality, i.e., high definition phenotyping

113    of naïve and memory subsets, remains a challenge.

114    Generally, T cells with a lower differentiation state such as naïve and stem cell or central memory

115    cells have been shown to provide superior anti-tumor potency, presumably due to their higher

116    potential to replicate, migrate, and engraft, leading to a long-term, durable response[18–21]. Likewise,

117    CD4 T cells are similarly important to anti-tumor potency due to their cytokine release properties

118    and ability to resist exhaustion[22,23]. Our group has developed a novel degradable microscaffold

119    (DMS)-based method using porous microcarriers functionalized with anti-CD3 and anti-CD28

120    mAbs for use in T cell expansion cultures. We showed that compared to commercially available

121    microbeads (Miltenyi), degradable microscaffolds (DMSs) generated a higher number of

122    migratory naïve ($T_N$) and central-memory ($T_{CM}$) ($CCR7^+CD62L^+$) T cells and $CD4^+$ T cells across

123    multiple donors[12]. We used this manufacturing process as an exemplar to develop an experimental-

124    computational AI-based tool to predict product quality from early process measurements. This

125    two-phase approach consists of (1) the optimization of process parameters through experimental

126    designs, and (2) the extraction of early predictive signatures of T cell quality by multi-omics

127    integration using regression models. This agnostic computational approach provides a platform to

128    discover early predictive CQAs and CPPs to ensure consistent product quality, that can be widely

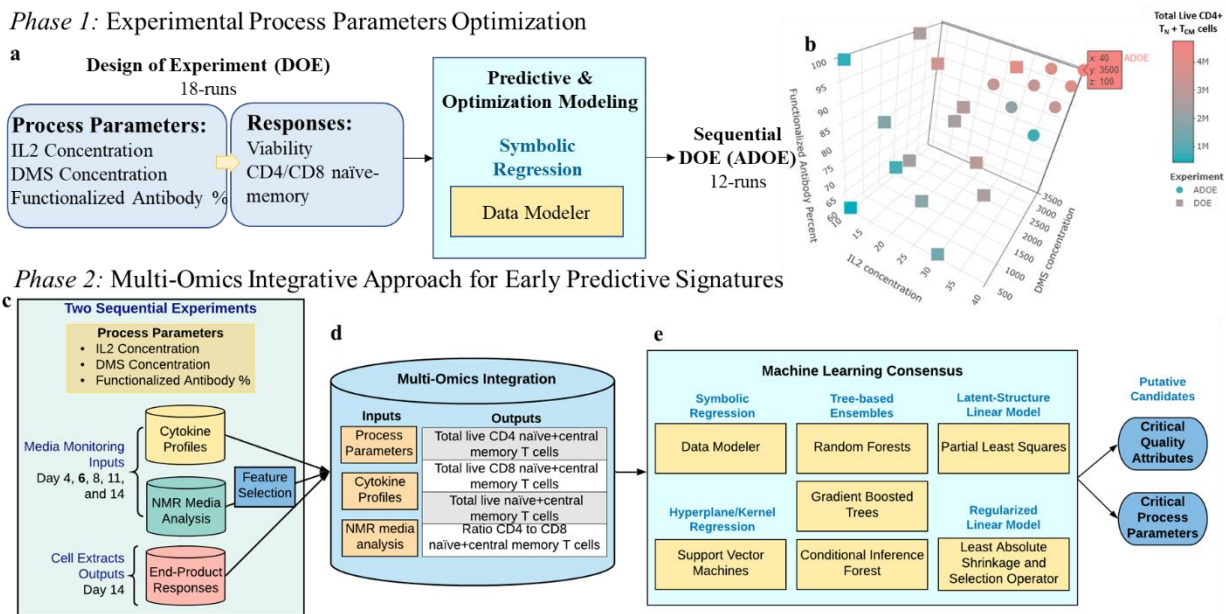129    applicable for other cellular therapies.

130

**Results**

132    I.    Overall multi-omics study design

133    T cells were expanded *ex vivo* for 14 days and 100 µL of supernatant media samples were collected

134    at days 4, 6, 8, 11, and 14 to measure cytokine profiles and perform NMR analysis. Endpoint

135    responses on DMS-based T cell extracts were measured for different combinations of DMS

136    parameters: IL2 concentration, DMS concentration, and functionalized antibody percent. Two

137    experimental regions were determined using a design-of-experiments (DOE) methodology to

138    maximize the yields of $CD62L^+CCR7^+$ cells (i.e. naïve and central memory T cells, $T_N+T_{CM}$) as a

139    function of these process parameters. The first DOE resulted in a randomized 18-run I-optimal

140    custom design where each DMS parameter was evaluated at three levels. To further optimize this

141    DOE in terms of total live $CD4^+$ $T_N+T_{CM}$ cells, a sequential adaptive design-of-experiment

142    (ADOE) was designed with 12 additional samples (Fig.1b). All 30 runs from both experiments

143    (DOE, ADOE) were molecularly characterized to model total live $T_N+T_{CM}$ (a) $CD4^+$, (b) $CD8^+$,

144    and (c) their ratio. The extraction of early predictive CPPs and CQAs for the expansion of $T_N+T_{CM}$

145    cells during *ex vivo* culture was performed in two phases: (1) optimization of process parameters,

146    and (2) integration of multi-omics for predictive modeling (Fig.1).



147

148    **Fig.1. Two-phase approach to extract early predictive CPPs and CQAs for CD4$^+$/CD8$^+$**

149    **T$_N$+T$_{CM}$ cells**. **a** DOE modeling and optimization of process parameters. **b** Experimental region

150    studied and optimized for total live CD4$^+$ T$_N$+T$_{CM}$ cells. **c** Overall study design (two experiments

151    varying process parameters while measuring multi-omics and T$_N$+T$_{CM}$ responses. e). **d** Integrative

152    multi-omics approach through **e** a machine learning consensus analysis to identify early predictive

153    CPPs and CQAs putative candidates for both total live CD4$^+$ and CD8$^+$ T$_N$+T$_{CM}$ cells.

154

155

156    II.    Optimization of T$_N$+T$_{CM}$ cells as a function of process parameters

157    Using symbolic regression (DataModeler software from Evolved Analytics LLC), we examined

158    the interactive effects of the DMS parameters on yield to simultaneously predict and optimize both

159    CD4$^+$ and CD8$^+$ T$_N$+T$_{CM}$. A model ensemble predicted 4.2 x 10$^6$ CD4$^+$ T$_N$+T$_{CM}$ cells at an

160    optimum setting of 30 U/μL IL2, 2500 carriers/μL, and 100% functionalized mAbs

161    (Supp.Fig.S1,S3,S4). This result was consistent with the observed maximum value of 4.0 x $10^6$,

162    highlighting that CD4$^+$ $T_N$+$T_{CM}$ yield was maximized at high levels of DMS parameters (Fig.1b).

163    In contrast, the predicted optimum yield for CD8$^+$ $T_N$+$T_{CM}$ was 1.9 x $10^7$ cells at a setting of 30

164    U/μL IL2, 600 carriers/μL, and 100% functionalized mAbs (Supp.Fig.S2,S3,S4). Although this

165    combination was not experimentally tested, the closest measured record (30 U/μL IL2, 500

166    carriers/μL, 100% functionalized mAbs) achieved the predicted maximum yield. Hence, the CD8$^+$

167    $T_N$+$T_{CM}$ yield was maximized at high IL2 concentration and functionalized mAbs percentage but

168    low DMS concentration.

169    The DOE analysis highlighted the potential for further optimization of total live CD4$^+$ $T_N$+$T_{CM}$

170    cells, as well as the potential to optimize the CD4$^+$ to CD8$^+$ $T_N$+$T_{CM}$ cells ratio, at DMS levels

171    greater than those originally evaluated (DOE). Therefore, to test and validate, a second adaptive

172    design of experiment (ADOE) was designed to maximize the total live CD4$^+$ $T_N$+$T_{CM}$ cells. We

173    expanded    the    parameter    range,    assessing    IL2    concentration>30    U/μL    and    DMS

174    concentration>2500 carriers/μL (Fig.1b). CD4$^+$ $T_N$+$T_{CM}$ and its ratio to CD8$^+$ $T_N$+ $T_{CM}$, 4.7 x $10^6$

175    cell and 0.49 respectively, were maximized when IL2 concentration (40 U/μL) and DMS

176    concentration (3500 carriers/μL) were maximized (Fig.1b;Supp.Table.S2;Supp.Fig.S1-S11).

177    Utilizing the ADOE dataset, new response ensembles were generated enabling more robust

178    prediction over the expanded parameter space (↑IL2 and ↑DMS concentrations).

179

180    III.    Multi-omic integrative analysis for early monitoring of T cell manufacturing

181    Due to the heterogeneity of the multivariate data collected and knowing that no single model

182    structure is perfect for all applications, we implemented an agnostic modeling approach to better

183    understand these $T_N+T_{CM}$ responses. To achieve this, a consensus analysis using seven machine

184    learning (ML) techniques, Random Forest (RF), Gradient Boosted Machine (GBM), Conditional

185    Inference Forest (CIF), Least Absolute Shrinkage and Selection Operator (LASSO), Partial Least-

186    Squares Regression (PLSR), Support Vector Machine (SVM), and DataModeler's Symbolic

187    Regression (SR), was implemented to molecularly characterize $T_N+T_{CM}$ cells and to extract

188    predictive features of quality early on their expansion process (Fig.1d-e).

189    SR models achieved the highest predictive performance ($R^2>93\%$) when using multi-omics

190    predictors for all endpoint responses (Table.1). SR achieved $R^2>98\%$ while GBM tree-based

191    ensembles showed leave-one-out cross-validated $R^2$ (LOO-$R^2$) >95% for $CD4^+$ and $CD4^+/CD8^+$

192    $T_N+T_{CM}$ responses. Similarly, LASSO, PLSR, and SVM methods showed consistent high LOO-

193    $R^2$, 92.9%, 99.7%, and 90.5%, respectively, to predict the $CD4^+/CD8^+$ $T_N+T_{CM}$. Yet, about 10%

194    reduction in LOO-$R^2$, 72.5%-81.7%, was observed for $CD4^+$ $T_N+T_{CM}$ with these three methods.

195    Lastly, SR and PLSR achieved $R^2>90\%$ while other ML methods exhibited exceedingly variable

196    LOO-$R^2$ (0.3%,RF-51.5%,LASSO) for $CD8^+$ $T_N+T_{CM}$ cells.

197    The top-performing technique, SR, showed that the median aggregated predictions for $CD4^+$ and

198    $CD8^+$ $T_N+T_{CM}$ cells increases when IL2 concentration, IL15, and IL2R increase while IL17a

199    decreases in conjunction with other features. These patterns combined with low values of DMS

200    concentration and GM_CSF uniquely characterized maximum $CD8^+$ $T_N+T_{CM}$. Meanwhile, higher

201    glycine but lower IL13 in combination with others showed maximum $CD4^+$ $T_N+T_{CM}$ predictions

202    (Fig.2).

203

**Table 1. LOO-R$^2$ prediction performance results for all ML models when evaluating process parameters, and features from cytokine and NMR media analysis at day 6 or day 4.**

| LOO-R2 Response/Predictors | ML | | | | | | |
|---|---|---|---|---|---|---|---|
| | SR | RF | GBM | CIF | LASSO | PLSR | SVM |
| **Ratio.of.CD4.to.CD8.TN+TCM.Cells** | | | | | | | |
| PP+N4 | **99%** | 86.8% | 96.3% | **84.5%** | 88.6% | 92.5% | 88.5% |
| PP+N6 | **99%** | 73.6% | 95.9% | 70.1% | 81.0% | 95.8% | 79.7% |
| PP+S6 | **99%** | **87.1%** | **99.9%** | 83.4% | 87.2% | 97.9% | 86.8% |
| PP+S6+N6 | **99%** | 85.5% | 95.3% | 83.4% | **92.9%** | **99.7%** | **90.5%** |
| **Total.live.CD4+.TN+TCM.cells** | | | | | | | |
| PP+N4 | 97% | 67.0% | 93.6% | 69.3% | 34.3% | 90.1% | 75.5% |
| PP+N6 | 96% | 45.9% | 92.6% | 51.2% | 42.8% | **92.1%** | **79.4%** |
| PP+S6 | **98%** | **71.4%** | **99.9%** | **75.0%** | **74.9%** | 80.0% | 75.5% |
| PP+S6+N6 | **98%** | 68.2% | 95.6% | 74.4% | 72.5% | 81.7% | 77.0% |
| **Total.live.CD8+.TN+TCM.cells** | | | | | | | |
| PP+N4 | 93% | 4.7% | **44.4%** | 9.2% | 1.2% | 65.1% | 9.1% |
| PP+N6 | 86% | 2.0% | 29.9% | **15.8%** | 28.5% | 63.3% | 30.6% |
| PP+S6 | **93%** | **7.8%** | 28.0% | 15.1% | **76.2%** | **98.4%** | **49.8%** |
| PP+S6+N6 | **93%** | 0.3% | 32.7% | 9.8% | 51.5% | 96.4% | 37.8% |

ML models prediction performance is measured as the leave-one-out cross-validated R$^2$ (LOO-R$^2$) while SR prediction performance is measured as R$^2$ of the ensemble prediction where the ensemble is composed of diverse models with complexity constrained. Predictors evaluated: (PP) Process parameters, (N) NMR, (S) Cytokines measured at day 4 or 6. **max R$^2$ within each ML method are shown in bold.**
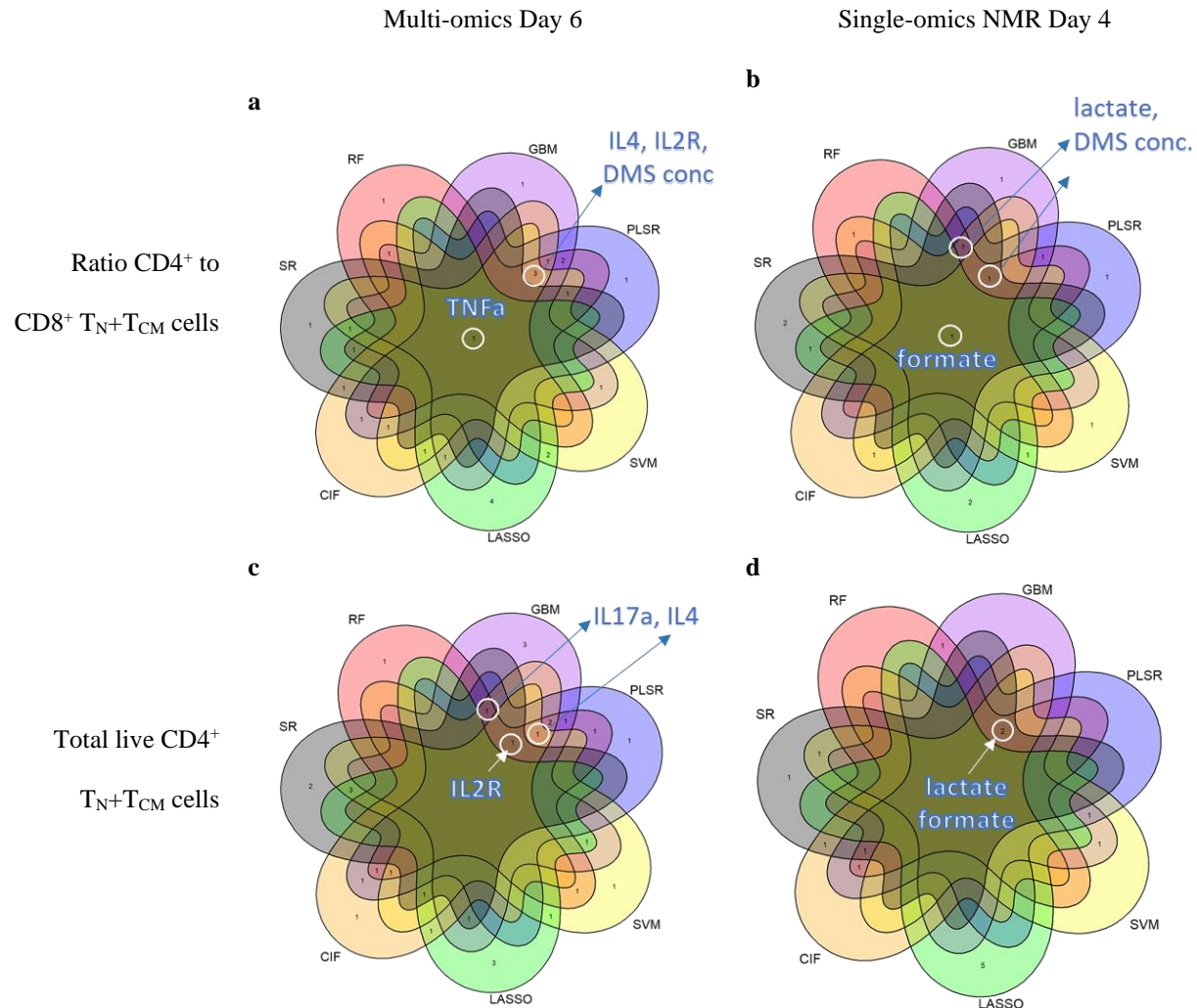
211

**Fig.2. Multi-omics culturing media prediction profiles at day 6 from DataModeler.** Prediction model profiles from day 6 culturing media monitoring where total live $CD4^+$ $T_N+T_{CM}$ is maximized.

Selecting CPPs and CQAs candidates consistently for T cell memory is desired. Here, TNFα was found in consensus across all seven ML methods for predicting $CD4^+/CD8^+$ $T_N+T_{CM}$ when considering features with the highest importance scores across models (Fig.3a;Methods). Other features, IL2R, IL4, IL17a, and DMS concentration, were commonly selected in $\geq 5$ ML methods (Fig.3a,c). Moreover, IL13 and IL15 were found predictive in combination with these using SR (Supp.Table.S4).

This integrative analysis of cytokine and NMR media analysis monitored at early stages of the T cell process provided highly predictive feature combinations of end-product quality. However, when translating a real-time monitoring strategy to a large-scale manufacturing process, measuring both cytokine and NMR features from media can be difficult and expensive. To be cost-efficient and translatable, we demonstrated that either cytokine profiles or NMR media analysis alone is sufficient to find predictive features without compromising prediction performance.

228



**Fig.3. ML model consensus of highly predictive for early monitoring of T cell manufacturing.**

ML models consensus for **a-b** ratio CD4$^+$ to CD8$^+$ T$_N$+T$_{CM}$ cells, and **c-d** total live CD4$^+$ T$_N$+T$_{CM}$ cells for both multi-omics modeling at day 6 and single-omics with NMR at day 4, respectively. Feature names are shown for consensus with 5 or more ML models at the highest-ranking standing (see Methods).

237    IV.    Cytokine media profiles for early prediction

238    ML models using solely media cytokine profiles at day 6 reached similar or higher $R^2$ than those

239    of the multi-omics models (CD4$^+$ $T_N$+$T_{CM}$: 71.4%-99.9%; CD4$^+$/CD8$^+$: 83.4%-99.7%). However,

240    CD8$^+$ $T_N$+$T_{CM}$ still had variable LOO-$R^2$, 7.8%-93%. Overall, higher cytokine media profiles

241    showed higher CD4$^+$ $T_N$+$T_{CM}$ and consequently its ratio with CD8$^+$ (Fig.4a). This behavior was

242    evident, even beyond day 6, for TNFα, IL2R, IL17a, and IL4 which were frequently selected as

243    predictive features across models (Fig.4b-c;Supp.Fig.S20). A more complex behavior was

244    detected for CD8$^+$ $T_N$+$T_{CM}$ which cannot be explained by cytokine secretion alone (Fig.4d).

245

246    V.    NMR media analysis for early prediction

247    Models using only NMR media intensities on day 6 revealed an $R^2$ decrease of 8.8% and 11.1%,

248    on average, compared with the multi-omics and cytokine models, respectively. Yet, SR, GBM,

249    and PLSR reached high LOO-$R^2$ (92.1%-99%), specifically for CD4$^+$/CD8$^+$ and CD4$^+$ $T_N$+$T_{CM}$.

250    Although good prediction was achieved with NMR media analysis on day 6, we obtain slightly

251    better predictions with NMR media analysis on day 4 (Table.1). From these models, formate,

252    lactate, DMS concentration were highly ranked to predict both, ratio CD4$^+$/CD8$^+$ and CD4$^+$

253    $T_N$+$T_{CM}$ (Fig.3b,d;Supp.Fig.19d). Some variable combinations also contained histidine, ethanol,

254    dimethylamine, branch chain amino acids (BCAAs), glucose, and glutamine (Supp.Table.S3).

255    Lower intensity values for BCAAs, dimethylamine, glucose, and glutamine displayed higher CD4$^+$

256    $T_N$+$T_{CM}$ cells across the different media monitoring times (Supp.Fig.S25). Inversely, higher

257    intensities of formate and lactate showed higher CD4$^+$ $T_N$+$T_{CM}$ and its ratio with CD8$^+$ consistently

258    across time (Fig.5a,b).

259 **a**



260
261 **b**



262
263 **c**



264
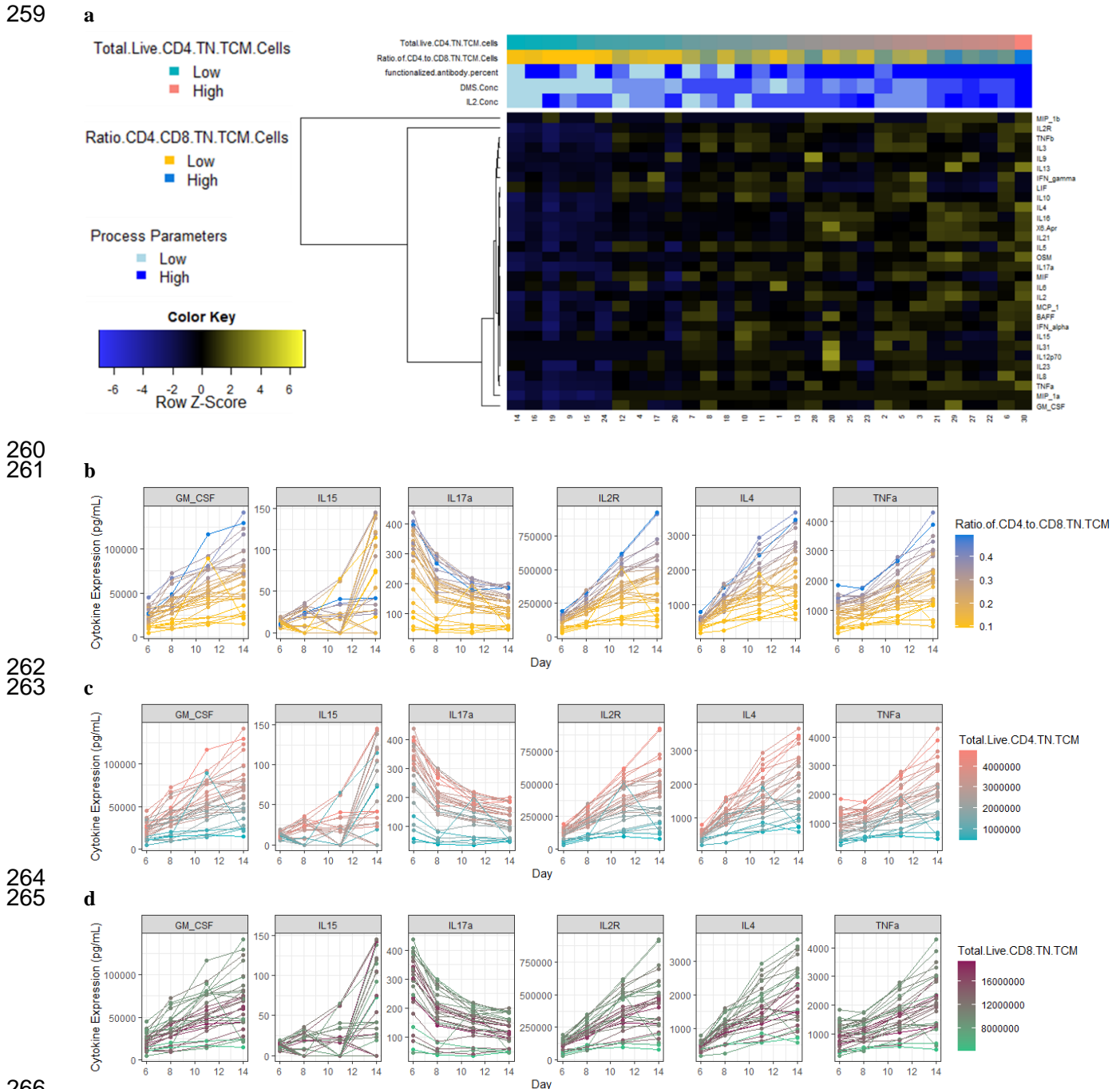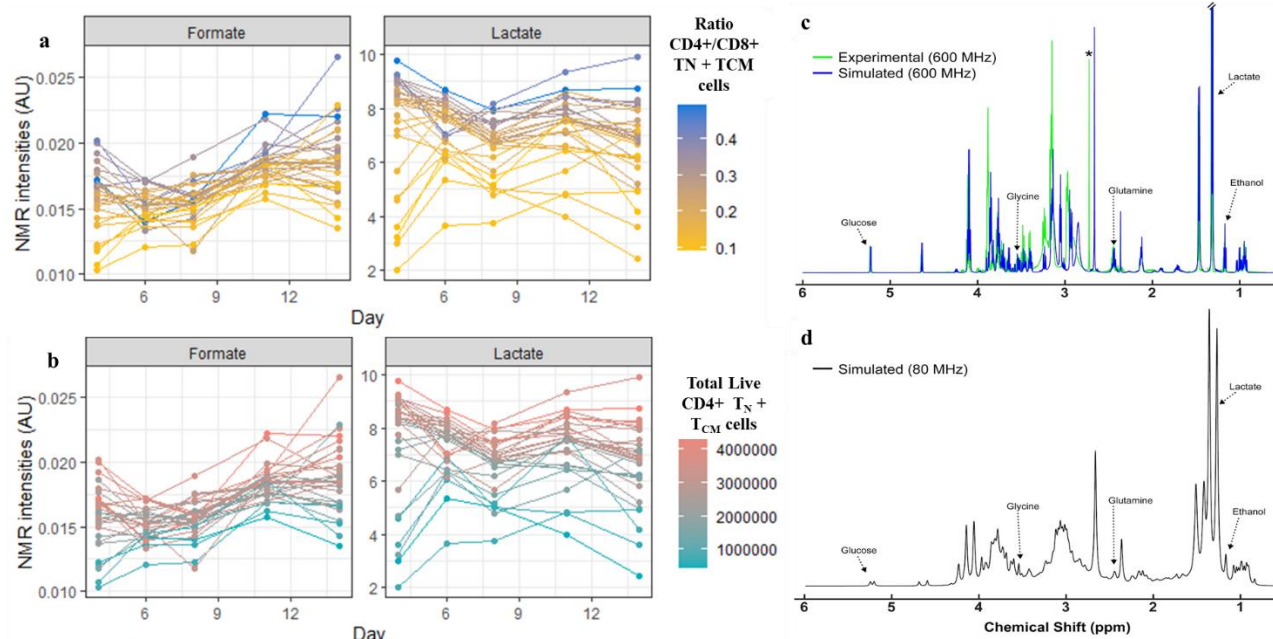265 **d**



266
267 **Fig.4. General characteristics of cytokine media profiles. a** Heatmap for cytokine profiles from

268 media samples on day 6. Expression in picograms/milliliter across time points for relevant

269 cytokine features for **b** ratio CD4+ to CD8+ $T_N$+$T_{CM}$ cells, **c** total live CD4$^+$ $T_N$+$T_{CM}$ cells, and **d**

270 total live CD8$^+$ $T_N$+$T_{CM}$ cells.

271

272

273



274

**Fig.5. Top-performing features NMR media analysis.** NMR intensities in arbitrary units (AU) across time points for **a** Ratio CD4$^+$/CD8$^+$ T$_N$+T$_{CM}$ cells, and **b** total live CD4$^+$ T$_N$+T$_{CM}$ cells. **c** Simulation of $^1$H NMR spectrum shows the potential to detect multiple predictive features at lower magnetic fields. Overlay of a pooled experimental spectrum of T-cell culture medium (green) and GISSMO[27,28] simulated spectrum (blue), composed of 19 compounds that reasonably approximate the experimental spectrum acquired at 600 MHz. *indicates an unknown feature of high intensity that was simulated with 2,3-dimethylamine (blue feature to right). Annotated features in the spectrum correspond to those identified as being highly predictive of output responses across computational methods. **d** GISSMO[27,28] simulated spectrum at 80 MHz, corresponding to a field strength of commercially available benchtop NMR systems.

285

286

**Discussion**

287

288 I.  Optimization of process parameters

289 CPPs modeling and understanding are critical to new product development and in cell therapy

290 development, it can have life-saving implications. The challenges for effective modeling grow

291 with the increasing complexity of processes due to high dimensionality, and the potential for

292 process interactions and nonlinear relationships. Another critical challenge is the limited amount

293 of available data, mostly small DOE datasets. SR has the necessary capabilities to resolve the

294 issues of process effects modeling and has been applied across multiple industries[12]. SR discovers

295 mathematical expressions that fit a given sample and differs from conventional regression

296 techniques in that a model structure is not defined *a priori*[13]. Hence, a key advantage of this

297 methodology is that transparent, human-interpretable models can be generated from small and

298 large datasets with no prior assumptions[14,15].

299 Since the model search process lets the data determine the model, diverse and competitive (e.g.,

300 accuracy, complexity) model structures are typically discovered. An ensemble of diverse models

301 can be formed where its constituent models will tend to agree when constrained by observed data

302 yet diverge in new regions. Collecting data in these regions helps to ensure that the target system

303 is accurately modeled, and its optimum is accurately located[14,15]. Exploiting these features allows

304 adaptive data collection and interactive modeling. Consequently, this adaptive-DOE approach is

305 useful in a variety of scenarios, including maximizing model validity for model-based decision

306 making, optimizing processing parameters to maximize target yields, and developing emulators

307 for online optimization and human understanding[14,15].

308

309 II.  Early predictive features

310    An in-depth characterization of potential DMS-based T-cell CQAs includes a list of cytokine and

311    NMR features from media samples that are crucial in many aspects of T cell fate decisions and

312    effector functions of immune cells. Cytokine features were observed to slightly improve prediction

313    and dominated the ranking of important features and variable combinations when modeling

314    together with NMR media analysis and process parameters (Fig.3b,d).

315    Predictive cytokine features such as TNFα, IL2R, IL4, IL17a, IL13, and IL15 were biologically

316    assessed in terms of their known functions and activities associated with T cells. T helper cells

317    secrete more cytokines than T cytotoxic cells, as per their main functions, and activated T cells

318    secrete more cytokines than resting T cells. It is possible that some cytokines simply reflect the

319    $CD4^+/CD8^+$ ratio and the activation degree by proxy proliferation. However, the exact ratio of

320    expected cytokine abundance is less clear and depends on the subtypes present, and thus

321    examination of each relevant cytokine is needed.

322    IL2R is secreted by activated T cells and binds to IL2, acting as a sink to dampen its effect on T

323    cells[16]. Since IL2R was much greater than IL2 in solution, this might reduce the overall effect of

324    IL2, which could be further investigated by blocking IL2R with an antibody. In T cells, TNF can

325    increase IL2R, proliferation, and cytokine production[18]. It may also induce apoptosis depending

326    on concentration and alter the $CD4^+$ to $CD8^+$ ratio[17]. Given that TNF has both a soluble and

327    membrane-bound form, this may either increase or decrease $CD4^+$ ratio and/or memory T cells

328    depending on the ratio of the membrane to soluble TNF[18]. Since only soluble TNF was measured,

329    membrane TNF is needed to understand its impact on both $CD4^+$ ratio and memory T cells.

330    Furthermore, IL13 is known to be critical for Th2 response and therefore could be secreted if there

331    are significant Th2 T cells already present in the starting population[19]. This cytokine has limited

332    signaling in T cells and is thought to be more of an effector than a differentiation cytokine[20]. It

333 might be emerging as relevant due to an initially large number of Th2 cells or because Th2 cells

334 were preferentially expanded; indeed, IL4, also found important, is the conical cytokine that

335 induces Th2 cell differentiation (Fig.3). The role of these cytokines could be investigated by

336 quantifying the Th1/2/17 subsets both in the starting population and longitudinally. Similar to

337 IL13, IL17 is an effector cytokine produced by Th17 cells[21] thus may reflect the number of Th17

338 subset of T cells. GM-CSF has been linked with activated T cells, specifically Th17 cells, but it is

339 not clear if this cytokine is inducing differential expansion of $CD8^+$ T cells or if it is simply a

340 covariate with another cytokine inducing this expansion[22]. Finally, IL15 has been shown to be

341 essential for memory signaling and effective in skewing CAR-T cells toward the Tscm phenotype

342 when using membrane-bound IL15Ra and IL15R[23]. Its high predictive behavior goes with its

343 ability to induce large numbers of memory T cells by functioning in an autocrine/paracrine manner

344 and could be explored by blocking either the cytokine or its receptor.

345 Moreover, many predictive metabolites found here are consistent with metabolic

346 activity associated with T cell activation and differentiation, yet it is not clear how the various

347 combinations of metabolites relate with each other in a heterogeneous cell population. Formate

348 and lactate were found to be highly predictive and observed to positively correlate with higher

349 values of total live $CD4^+$ $T_N+T_{CM}$ cells (Fig.5a-b;Supp.Fig.28-S30,S38). Formate is a byproduct

350 of the one-carbon cycle implicated in promoting T cell activation[24]. Importantly, this cycle occurs

351 between the cytosol and mitochondria of cells and formate excreted[25]. Mitochondrial biogenesis

352 and function are shown necessary for memory cell persistence[26,27]. Therefore, increased formate

353 in media could be an indicator of one-carbon metabolism and mitochondrial activity in the culture.

354 In addition to formate, lactate was found as a putative CQA of $T_N+T_{CM}$. Lactate is the end-product

355 of aerobic glycolysis, characteristic of highly proliferating cells and activated T cells[28,29]. Glucose

356 import and glycolytic genes are immediately upregulated in response to T cell stimulation, and

357 thus generation of lactate. At earlier time-points, this abundance suggests a more robust induction

358 of glycolysis and higher overall T cell proliferation. Interestingly, our models indicate that higher

359 lactate predicts higher CD4$^+$, both in total and in proportion to CD8$^+$, seemingly contrary to

360 previous studies showing that CD8$^+$ T cells rely more on glycolysis for proliferation following

361 activation[30]. It may be that glycolytic cells dominate in the culture at the early time points used for

362 prediction, and higher lactate reflects more cells.

363 Ethanol patterns are difficult to interpret since its production in mammalian cells is still poorly

364 understood[31]. Fresh media analysis indicates ethanol presence in the media used, possibly utilized

365 as a carrier solvent for certain formula components. However, this does not explain the high

366 variability and trend of ethanol abundance across time (Supp.Fig.S25-S27). As a volatile chemical,

367 variation could be introduced by sample handling throughout the analysis process. Nonetheless, it

368 is also possible that ethanol excreted into media over time, impacting processes regulating redox

369 and reactive oxygen species which have previously been shown to be crucial in T cell signaling

370 and differentiation[32].

371 Metabolites that consistently decreased over time are consistent with the primary carbon source

372 (glucose) and essential amino acids (BCAA, histidine) that must be continually consumed by

373 proliferating cells. Moreover, the inclusion of glutamine in our predictive models also suggests the

374 importance of other carbon sources for certain T cell subpopulations. Glutamine can be used for

375 oxidative energy metabolism in T cells without the need for glycolysis[30]. Overall, these results are

376 consistent with existing literature that show different T cell subtypes require different relative

377 levels of glycolytic and oxidative energy metabolism to sustain the biosynthetic and signaling

378 needs of their respective phenotypes[33,34]. It is worth noting that the trends of metabolite abundance

379   here are potentially confounded by the partial replacement of media that occurred periodically

380   during expansion (Methods), thus likely diluting some metabolic byproducts (i.e. formate, lactate)

381   and elevating depleted precursors (i.e. glucose, amino acids). More definitive conclusions of

382   metabolic activity across the expanding cell population can be addressed by a closed system,

383   ideally with on-line process sensors and controls for formate, lactate, along with ethanol and

384   glucose.

385   III.    Monitoring of T-cell manufacturing with benchtop NMR systems

386   We demonstrated the ability to identify predictive markers using high-magnetic field NMR

387   spectrometers. However, these are expensive, require a significant amount of resources to house

388   and maintain, and would be the unlikely option for routine monitoring in industrial cell-

389   manufacturing. Another common method, liquid chromatography (LC) coupled to mass

390   spectrometry, has the advantage of a relatively smaller footprint and less upfront cost but it has

391   other drawbacks such as destruction of the sample and difficulty with components in culture media

392   that damage LC columns without extraction. Nevertheless, methods like continuous closed-loop

393   sampling are being developed to address this and might be readily available in the future[35].

394   Recently, permanent magnet-based NMR spectrometers (benchtop-size) have become available at

395   a lower cost. Many of these are readily configured for flow-through reaction monitoring, which

396   can be leveraged in a closed-cell manufacturing process. To explore the feasibility of such system,

397   we utilized a spectral simulation to evaluate if putative CQAs identified here could theoretically

398   be observed and quantified at a magnetic field strength of 80 MHz (benchtop systems). First, the

399   experimental data acquired at 600 MHz was approximated by creating a simulated mixture of

400   identified metabolites (Fig.5c) and then simulated at 80 MHz (Fig.5d). While the spectral

401   resolution is significantly reduced compared to a spectrum at high-field, there are still numerous

402   features that can be attributed to unique metabolites, including those identified as highly predictive

403   (Fig.5c,d). Although this is promising, there will be challenges to acquiring high-quality data in a

404   closed bioreactor system, i.e. cells/DMS-particles in suspension, media formulation dictated by

405   spectral complexity/overlap, and accurate quantitation of features with high overlap from other

406   signals. However, a dedicated benchtop NMR coupled to a bioreactor could provide a simple

407   system for real-time monitoring of CQAs.

408   Henceforth, this two-phase approach enabled in-depth characterization and identification of

409   potential CQAs and CPPs for T cells. More sampling is needed to explore aspects like donor-to-

410   donor variability, when available it can be incorporated into this workflow which will be enriched

411   due to its data-driven iterative design that fine-tunes model parameters as more data fits back into

412   it. Providing a powerful framework to optimize a complex experimental space during the cell-

413   manufacturing process, and to facilitate the identification of CPPs and early predictive CQAs from

414   multi-omics, that can be used broadly in the cell therapy and regenerative medicine field to

415   accurately predict end-of-manufacturing quality at early stages.

416

417   **Methods**

418   I.   Overall multi-omics study design and development: More details

419   The first DOE resulted in a randomized 18-run I-optimal custom design where each DMS

420   parameter was evaluated at three levels: IL2 concentration (10, 20, and 30 U/µL), DMS

421   concentration (500, 1500, 2500 carrier/µL), and functionalized antibody percent (60%, 80%,

422   100%). These 18 runs consisted of 14 unique parameter combinations where 4 of them were

423   replicated twice to assess prediction error. Process parameters for the ADOE were evaluated at

424   multiple levels: IL2 concentration (30, 35, and 40 U/µL), DMS concentration (500, 1000, 1500,

425     2000, 2500, 3000, 3500 carrier/µL), and functionalized antibody percent (100%) as depicted in

426     Fig.1b. To further optimize the initial region explored (DOE) in terms of total live CD4$^+$ $T_N$+$T_{CM}$

427     cells, a sequential adaptive design-of-experiment (ADOE) was designed with 10 unique parameter

428     combinations, two of these replicated twice for a total of 12 additional samples (Fig.1b). The fusion

429     of cytokine and NMR profiles from media to model these responses included 30 cytokines from a

430     custom Thermo Fisher ProcartaPlex Luminex kit and 20 NMR features. These 20 spectral features

431     from NMR media analysis were selected out of approximately 250 peaks through the

432     implementation of a variance-based feature selection approach and some manual inspection steps.

433

434     II.     Microcarrier fabrication

435     Degradable microscaffolds were fabricated as previously described[36]. Briefly, gelatin

436     microcarriers (CuS, GE Healthcare DG-2001-OO) were suspended at 20 mg/mL in 1X phosphate-

437     buffered saline (PBS). Sulfo-NHS-biotin (SNB) (Thermo Fisher 21217 or Apex Bio A8001) was

438     dissolved at 10 µM in ultrapure water and 7.5 µL SNB/mL PBS was added to carrier suspension

439     and allowed to react for 60 min. After washing the carriers three times in PBS, 40 µg/mL

440     streptavidin (Jackson Immunoresearch 016-000-114) was added and allowed to react for 60 min.

441     Biotinylated mAbs against human CD3 and CD28 were combined in a 1:1 mass ratio and added

442     to the carriers at 2 µg mAbs/mg carriers. To vary the surface concentration of the antibodies, the

443     anti-CD3/anti-CD28 mAb mixture was further combined with a biotinylated isotype control to

444     reduce the overall fraction of targeted mAbs. mAbs were allowed to bind to the carriers for 60

445     min. All mAbs were low endotoxin azide-free (Biolegend custom, LEAF specification). Fully

446     functionalized DMSs were washed in sterile PBS and washed once again in the cell culture media

447   to be used for the T cell expansion. The surface concentration of the antibodies was quantified as

448   previously described using a bicinchoninic acid assay (BCA) kit (Thermo Fisher 23227)[36].

449

450   III.   T cell culture (including sample collection)

451   Cryopreserved primary human T cells were obtained as sorted CD3 subpopulations (Astarte

452   Biotech). T cells were activated by adding DMSs (amount specified by the DOE) at day 0 of culture

453   immediately after thaw. DMSs were not added or removed during the culture and had antibodies

454   that were conjugated in proportions specified by the DOE. Initial cell density was $2.0 * 10^6$ cells/mL

455   in a 96 well plate with 300 µL volume. Media was serum-free TexMACS (Miltentyi Biotech 170-

456   076-307) supplemented with recombinant human IL2 in concentrations specified by the DOE

457   (Peprotech 200-02). Cell cultures were expanded for 14 days as counted from the time of initial

458   seeding and activation. Cell counts and viability were assessed using acridine orange/propidium

459   iodide (AO/PI) and a Countess Automated Cell Counter (Thermo Fisher). Media was added to

460   cultures every 2 days to 3 days in a 3:1 ratio (new volume: old volume) or based on a 300 mg/dL

461   glucose threshold. The ADOE was done using the same feeding schedule as the initial DOE to

462   maintain consistency for validation. Media glucose was measured using a ChemGlass glucometer

463   to confirm cell growth and activation.

464

465   IV.   Flow cytometry

466   At the end of culture, at least 1e5 T cells from each run were washed with PBS once, resuspended

467   in PBS, and stained with Zombie UV (Biolegend, 423107) for 30 minutes at room temperature in

468   the dark at a 1:1000 dilution. Cells were spun and resuspended in FACS buffer (1X PBS, 2%

469   bovine serum albumin, 5 mM EDTA) and were stained with antibodies according to **Table M1** for

470 60 minutes in the dark at 4C. Cells were then resuspended in fresh FACS buffer, after which they

471 were run on a BD LSR ortessa. All stained was performed in a 96 well v-bottom plate.

472 **Table M1: Flow cytometry antibodies**

| Antigen | Fluorophore | Vendor | Cat Number |
|---|---|---|---|
| CD3 | APC-Fire | Biolegend | 34839 |
| CD4 | PerCP-Cy5.5 | BD | 561438 |
| CCR7 | AF647 | BD | 561438 |
| CD62L | PE | BD | 341012 |

473

474 V. Cytokine measurements

475 Cytokines were measured using a custom ProcartaPlex Luminex kit (Thermo Fisher). The assay

476 was performed using media samples taken at various time points throughout the T cell culture

477 according to the manufacturer's instructions with modifications to half the reagent requirements.

478 Briefly, an 8 point standard curve was created with all included standards. 25 μL magnetic beads

479 were added to all required wells and washed three times. 25 μL of each standard or sample was

480 added to the wells and the plate was sealed and spun at 850 rpm for 120 minutes followed by three

481 washes. 12.5 μL detection antibody was added followed by sealing the plate and spinning for 60

482 minutes at 850 rpm and three washes. 25 μL streptavidin PE was added followed by the same spin

483 and wash steps. 120 μL of reading buffer was added to the plate, the plate was analyzed on a

484 BioPlex 200 (BioRad). Any samples that were majority over-range (denoted as "OOR >" in the

485 output spreadsheet) were deemed too concentrated at run at 1/10th their original concentration to

486 put them within range. All samples were run without technical replicates.

487 Luminex data was preprocessed using R for inclusion in the analysis pipeline as follows. Any

488 cytokine level that was over-range ("OOR >" in output) was set to the maximum value of the

489    standard curve for that cytokine. Any value that was under-range ("OOR <" in output spreadsheet)

490    was set to zero. All values that were extrapolated from the standard curve were left unchanged.

491

492    VI.    NMR metabolomics

493    A.    Sample preparation

494    50 μL of media was collected from each culture at each time point (before media exchange, if

495    applicable), flash-frozen in liquid nitrogen, and stored at -80°C. Samples were shipped to CCRC

496    on dry ice for NMR analysis. Run order of samples was randomized. Samples were prepared in

497    two batches for each rack of NMR samples to be run. For each rack, samples were pulled and

498    sorted on dry ice, then thawed at 4°C for 1 hour. Samples were then centrifuged at 2,990 x g at

499    4°C for 20 minutes to pellet any cells or debris that may have been collected with the media. 5

500    μL of 100/3 mM DSS-D6 in deuterium oxide (Cambridge Isotope Laboratories) were added to

501    1.7 mm NMR tubes (Bruker BioSpin), followed by 45 μL of media from each sample that was

502    added and mixed, for a final volume of 50 μL in each tube. Samples were prepared on ice and in

503    predetermined, randomized order. The remaining volume from each sample in the rack (~4 μL)

504    was combined to create an internal pool. This material was used for internal controls within each

505    rack as well as metabolite annotation.

506    B.    Data collection

507    NMR spectra were collected on a Bruker Avance III HD spectrometer at 600 MHz using a 5-mm

508    TXI cryogenic probe and TopSpin software (Bruker BioSpin). One-dimensional spectra were

509    collected on all samples using the noesypr1d pulse sequence under automation using ICON NMR

510    software. Two-dimensional HSQC and TOCSY spectra were collected on internal pooled control

511    samples for metabolite annotation.

512    C.    Data processing

513    One-dimensional spectra were manually phased and baseline corrected in TopSpin. Two-

514    dimensional spectra were processed in NMRpipe[37]. One dimensional spectra were referenced,

515    water/end regions removed, and normalized with the PQN algorithm[38] using an in-house

516    MATLAB              (The             MathWorks,             Inc.)             toolbox

517    (https://github.com/artedison/Edison_Lab_Shared_Metabolomics_UGA).

518    D.    Feature selection

519    To reduce the total number of spectral features from approximately 250 peaks and enrich for those

520    that would be most useful for statistical modeling, a variance-based feature selection was

521    performed within MATLAB. For each digitized point on the spectrum, the variance was

522    calculated across all experimental samples and plotted. Clearly-resolved features corresponding

523    to peaks in the variance spectrum were manually binned and integrated to obtain quantitative

524    feature intensities across all samples (Supp.Fig.S24). In addition to highly variable features,

525    several other clearly resolved and easily identifiable features were selected (glucose, BCAA

526    region, etc). Some features were later discovered to belong to the same metabolite but were

527    included in further analysis.

528    E.    Metabolite annotation

529    Two-dimensional spectra collected on pooled samples were uploaded to COLMARm web

530    server[10], where HSQC peaks were automatically matched to database peaks. HSQC matches were

531    manually reviewed with additional 2D and proton spectra to confirm the match. Annotations were

532    assigned a confidence score based upon the levels of spectral data supporting the match as

533    previously described[11]. Annotated metabolites were matched to previously selected features used

534    for statistical analysis.

535    F.    Low-field spectrum simulation

536    Using the list of annotated metabolites obtained above, an approximation of a representative

537    experimental spectrum was generated using the GISSMO mixture simulation tool.[39,40] With the

538    simulated mixture of compounds, generated at 600 MHz to match the experimental data, a new

539    simulation was generated at 80 MHz to match the field strength of commercially available

540    benchtop NMR spectrometers. The GISSMO tool allows visualization of signals contributed from

541    each individual compound as well as the mixture, which allows annotation of features in the

542    mixture belonging to specific compounds.

543    G.  Unknown identification

544    Several low abundance features selected for analysis did not have database matches and were not

545    annotated. Statistical total correlation spectroscopy[41] suggested that some of these unknown

546    features belonged to the same molecules (not shown). Additional multidimensional NMR

547    experiments will be required to determine their identity.

548

549    VII.    Machine learning techniques & statistical analysis

550    A.  Machine learning modeling

551    Seven machine learning (ML) techniques were implemented to predict three responses related to

552    the memory phenotype of the cultured T cells under different process parameters conditions (i.e.

553    Total Live CD4+ $T_N$ and $T_{CM}$, Total Live CD8+ $T_N$+$T_{CM}$, and Ratio CD4+/CD8+ $T_N$+$T_{CM}$). The

554    ML methods executed were Random Forest (RF), Gradient Boosted Machine (GBM), Conditional

555 Inference Forest (CIF), Least Absolute Shrinkage and Selection Operator (LASSO), Partial Least-

556 Squares Regression (PLSR), Support Vector Machine (SVM), and DataModeler's Symbolic

557 Regression (SR). Primarily, SR models were used to optimize process parameter values based on

558 $T_N+T_{CM}$ phenotype and to extract early predictive variable combinations from the multi-omics

559 experiments. Furthermore, all regression methods were executed, and the high-performing models

560 were used to perform a consensus analysis of the important variables to extract potential critical

561 quality attributes and critical process parameters predictive of T-cell potency, safety, and

562 consistency at the early stages of the manufacturing process.

563 Symbolic regression (SR) was done using Evolved Analytics' DataModeler software (Evolved

564 Analytics LLC, Midland, MI). DataModeler utilizes genetic programming to evolve symbolic

565 regression models (both linear and non-linear) rewarding simplicity and accuracy. Using the

566 selection criteria of highest accuracy ($R^2$>90% or noise-power) and lowest complexity, the top-

567 performing models were identified. Driving variables, variable combinations, and model

568 dimensionality tables were generated. The top-performing variable combinations were used to

569 generate model ensembles. In this analysis, DataModeler's *SymbolicRegression* function was used

570 to develop explicit algebraic (linear and nonlinear) models. The fittest models were analyzed to

571 identify the dominant variables using the *VariablePresence* function, the dominant variable

572 combinations using the *VariableCombinations* function, and the model dimensionality (number of

573 unique variables) using the *ModelDimensionality* function. *CreateModelEnsemble* was used to

574 define trustable model ensembles using selected variable combinations and these were summarized

575 (model expressions, model phenotype, model tree plot, ensemble quality, model quality, variable

576 presence map, ANOVA tables, model prediction plot, exportable model forms) using the

577 *ModelSummaryTable* function. Ensemble prediction and residual performance were respectively

578    assessed via the *EnsemblePredictionPlot* and *EnsembleResidualPlot* subroutines. Model maxima

579    (*ModelMaximum* function) and model minima (*ModelMinimum* function) were calculated and

580    displayed using the *ResponsePlotExplorer* function. Trade-off performance of multiple responses

581    was explored using the *MultiTargetResponseExplorer* and *ResponseComparisonExplorer* with

582    additional insights derived from the *ResponseContourPlotExplorer*. Graphics and tables were

583    generated by DataModeler. These model ensembles were used to identify predicted response

584    values, potential optima in the responses, and regions of parameter values where the predictions

585    diverge the most.

586    Non-parametric tree-based ensembles were done through the *randomForest, gbm,* and *cforest*

587    regression functions in R, for random forest, gradient boosted trees, and conditional inference

588    forest models, respectively. Both random forest and conditional inference forest construct multiple

589    decision trees in parallel, by randomly choosing a subset of features at each decision tree split, in

590    the training stage. Random forest individual decision trees are split using the Gini Index, while

591    conditional inference forest uses a statistical significance test procedure to select the variables at

592    each split, reducing correlation bias. In contrast, gradient boosted trees construct regression trees

593    in series through an iterative procedure that adapts over the training set. This model learns from

594    the mistakes of previous regression trees in an iterative fashion to correct errors from its precursors'

595    trees (i.e. minimize mean squared errors). Prediction performance was evaluated using leave-one-

596    out cross-validation (LOO)-$R^2$ and permutation-based variable importance scores assessing %

597    increase of mean squared errors (MSE), relative influence based on the increase of prediction error,

598    coefficient values for RF, GBM, and CID, respectively. Partial least squares regression was

599    executed using the *plsr* function from the *pls* package in R while LASSO regression was performed

600    using the *cv.glmnet* R package, both using leave-one-out cross-validation. Finally, the *kernlab* R

601    package was used to construct the Support Vector Machine regression models.

602    Parameter tuning was done for all models in a grid search manner using the *train* function from

603    the *caret* R package using LOO-$R^2$ as the optimization criteria. Specifically, the number of features

604    randomly sampled as candidates at each split (mtry) and the number of trees to grow (ntree) were

605    tuned parameters for random forest and conditional inference forest. In particular, minimum sum

606    of weights in a node to be considered for splitting and the minimum sum of weights in a terminal

607    node were manually tuned for building the CIF models. Moreover, GBM parameters such as the

608    number of trees to grow, maximum depth of each tree, learning rate, and the minimal number of

609    observations at the terminal node, were tuned for optimum LOO-$R^2$ performance as well. For

610    PLSR, the optimal number of components to be used in the model was assessed based on the

611    standard error of the cross-validation residuals using the function *selectNcomp* from the *pls*

612    package. Moreover, LASSO regression was performed using the *cv.glmnet* package with *alpha* =

613    1. The best lambda for each response was chosen using the minimum error criteria. Lastly, a fixed

614    linear kernel (i.e. svmLinear) was used to build the SVM regression models evaluating the cost

615    parameter value with best LOO-$R^2$. Prediction performance was measured for all models using the

616    final model with LOO-$R^2$ tuned parameters. **Table M2** shows the parameter values evaluated per

617    model at the final stages of results reporting.

618

619    **Table M2: ML parameter values evaluated**

| ML Model | Tuned Parameter Values |
|---|---|
| RF | `ntree=c(500,1000,1500,2000,2500)`<br>`mtry=all possibilities` |
| GBM | `interaction.depth=c(1:4)`<br>`n.trees = (1:20)*10`<br>`shrinkage=c(0.1,0.01, 0.02)`<br>`n.minobsinnode=c(2:6)`<br>`bag.fraction=0.5` |

| CIF | `mtry=all possibilities`<br>`ntree*=100`<br>`minsplit* = 6`<br>`minbucket* = 3` |
|---|---|
| LASSO | `alpha=1`<br>`lambda=seq(0.001,0.05,by = 0.001)` |
| PLSR | `ncomp = 1:15` |
| SVM | `svmLinear`<br>`cost=seq(0.05,2,.05)` |
| | *other values besides the ones shown were optimized manually |

620

621 B. Consensus analysis

622 Consensus analysis of the relevant variables extracted from each machine learning model was done

623 to identify consistent predictive features of quality at the early stages of manufacturing. First

624 importance scores for all features were measured across all ML models using *varImp* with *caret* R

625 package except for scores for SVM which *rminer* R package was used. These importance scores

626 were percent increase in mean squared error (MSE), relative importance through average increase

627 in prediction error when a given predictor is permuted, permuted coefficients values, absolute

628 coefficient values, weighted sum of absolute coefficients values, and relative importance from

629 sensitivity analysis determined for RF, GBM, CIF, LASSO, PLSR, and SVM, respectively. Using

630 these scores, key predictive variables were selected if their importance scores were within the $80^{th}$

631 percentile ranking for the following ML methods: RF, GBM, CIF, LASSO, PLSR, SVM while for

632 SR variables present in >30% of the top-performing SR models from DataModeler (R2≥ 90%,

633 Complexity ≤ 100) were chosen to investigate consensus except for NMR media models at day 4

634 which considered a combination of the top-performing results of models excluding lactate ppms,

635 and included those variables which were in > 40% of the best performing models. Only variables

636 with those high percentile scoring values were evaluated in terms of their logical relation

637 (intersection across ML models) and depicted using a Venn diagram from the *venn* R package.

638

639 **Data availability**

640 The pre-processed set of the data used in this work is available in Supplementary Methods. All

641 NMR data are available at the Metabolomics Workbench[42] with DOI:

642 http://dx.doi.org/10.21228/M8F982.

643 **Code availability**

644 Machine learning implementation codes used in this work are available at GitHub

645 (https://github.com/wandaliz/CMaT_TCell_MachineLearning/). DataModeler information can be

646 requested at http://www.evolved-analytics.com/.

647

**References**

1. Fesnak, A. D., June, C. H. & Levine, B. L. Engineered T cells: the promise and challenges of cancer immunotherapy. *Nat. Rev. Cancer* **16**, 566–81 (2016).

2. Rosenberg, S. A. & Restifo, N. P. Adoptive cell transfer as personalized immunotherapy for human cancer. *Science* **348**, 62–8 (2015).

3. Dwarshuis, N. J., Parratt, K., Santiago-Miranda, A. & Roy, K. Cells as advanced therapeutics: State-of-the-art, challenges, and opportunities in large scale biomanufacturing of high-quality cells for adoptive immunotherapies. *Adv. Drug Deliv. Rev.* **114**, 222–239 (2017).

4. Roddie, C., O'Reilly, M., Pinto, J. D. A., Vispute, K. & Lowdell, M. Manufacturing chimeric antigen receptor T cells: issues and challenges. *Cytotherapy* (2019) doi:10.1016/j.jcyt.2018.11.009.

5. Roh, K.-H., Nerem, R. M. & Roy, K. Biomanufacturing of Therapeutic Cells: State of the Art, Current Challenges, and Future Perspectives. *Annu. Rev. Chem. Biomol. Eng.* **7**, 455–478 (2016).

6. Carmen, J., Burger, S. R., McCaman, M. & Rowley, J. A. Developing assays to address identity, potency, purity and safety: cell characterization in cell therapy process development. *Regen. Med.* **7**, 85–100 (2012).

7. Simon, C. G., Lin-Gibson, S., Elliott, J. T., Sarkar, S. & Plant, A. L. Strategies for Achieving Measurement Assurance for Cell Therapy Products. *Stem Cells Transl. Med.* **5**, 705–708 (2016).

8. Campbell, A. *et al.* Concise Review: Process Development Considerations for Cell Therapy. *Stem Cells Transl. Med.* **4**, 1155–1163 (2015).

9. Lipsitz, Y. Y., Timmins, N. E. & Zandstra, P. W. Quality cell therapy manufacturing by design. *Nat. Biotechnol.* **34**, 393–400 (2016).

673    10. Better, M., Chiruvolu, V. & Sabatino, M. Overcoming Challenges for Engineered Autologous

674        T Cell Therapies. *Cell Gene Ther. Insights* **4**, 173–186.

675    11. Tyagarajan, S., Spencer, T. & Smith, J. Optimizing CAR-T Cell Manufacturing Processes

676        during Pivotal Clinical Trials. *Mol. Ther. - Methods Clin. Dev.* **16**, 136–144 (2020).

677    12. Kordon, A. K. & Ching-Tai Lue. Symbolic regression modeling of blown film process effects.

678        in *Proceedings of the 2004 Congress on Evolutionary Computation (IEEE Cat.*

679        *No.04TH8753)* vol. 1 561-568 Vol.1 (2004).

680    13. Koza, J. R. & Koza, J. R. *Genetic Programming: On the Programming of Computers by*

681        *Means of Natural Selection.* (MIT Press, 1992).

682    14. Kotanchek, M. Adaptive Design-of-Experiments. 15.

683    15. Kotanchek, M., Smits, G. & Vladislavleva, E. Exploiting Trustable Models via Pareto GP for

684        Targeted Data Collection. in *Genetic Programming Theory and Practice VI* 1–18 (Springer

685        US, 2009). doi:10.1007/978-0-387-87623-8_10.

686    16. Witkowska, A. M. On the Role of sIL-2R Measurements in Rheumatoid Arthritis and

687        Cancers. *Mediators of Inflammation* vol. 2005 121–130

688        https://www.hindawi.com/journals/mi/2005/742316/ (2005).

689    17. Vudattu, N. K. *et al.* Reverse signalling of membrane-integrated tumour necrosis factor

690        differentially regulates alloresponses of CD4+ and CD8+ T cells against human

691        microvascular endothelial cells. *Immunology* **115**, 536–543 (2005).

692    18. Mehta, A. K., Gracias, D. T. & Croft, M. TNF activity and T cells. *Cytokine* **101**, 14–18

693        (2018).

694    19. Wong, F. S. Stimulating IL-13 Receptors on T cells: A New Pathway for Tolerance Induction

695        in Diabetes?-. *Diabetes* **60**, 1657–1659 (2011).

696    20. Junttila, I. S. Tuning the Cytokine Responses: An Update on Interleukin (IL)-4 and IL-13

697        Receptor Complexes. *Front. Immunol.* **9**, (2018).

698    21. Amatya, N., Garg, A. V. & Gaffen, S. L. IL-17 Signaling: The Yin and the Yang. *Trends*

699         *Immunol.* **38**, 310–322 (2017).

700    22. Becher, B., Tugues, S. & Greter, M. GM-CSF: From Growth Factor to Central Mediator of

701         Tissue Inflammation. *Immunity* **45**, 963–973 (2016).

702    23. Hurton, L. V. *et al.* Tethered IL-15 augments antitumor activity and promotes a stem-cell

703         memory subset in tumor-specific T cells. *Proc. Natl. Acad. Sci. U. S. A.* **113**, E7788–E7797

704         (2016).

705    24. Ron-Harel, N. *et al.* Mitochondrial biogenesis and proteome remodeling promotes one

706         carbon metabolism for T cell activation. *Cell Metab.* **24**, 104–117 (2016).

707    25. Pietzke, M., Meiser, J. & Vazquez, A. Formate metabolism in health and disease. *Mol.*

708         *Metab.* **33**, 23–37 (2020).

709    26. van der Windt, G. J. W. *et al.* Mitochondrial Respiratory Capacity Is A Critical Regulator Of

710         CD8+ T Cell Memory Development. *Immunity* **36**, 68–78 (2012).

711    27. Vardhana, S. A. *et al.* Impaired mitochondrial oxidative phosphorylation limits the self-

712         renewal of T cells exposed to persistent antigen. *Nat. Immunol.* 1–12 (2020)

713         doi:10.1038/s41590-020-0725-2.

714    28. Lunt, S. Y. & Vander Heiden, M. G. Aerobic Glycolysis: Meeting the Metabolic

715         Requirements of Cell Proliferation. *Annu. Rev. Cell Dev. Biol.* **27**, 441–464 (2011).

716    29. Chang, C.-H. *et al.* Posttranscriptional Control of T Cell Effector Function by Aerobic

717         Glycolysis. *Cell* **153**, 1239–1251 (2013).

718    30. Cao, Y., Rathmell, J. C. & Macintyre, A. N. Metabolic Reprogramming towards Aerobic

719         Glycolysis Correlates with Greater Proliferative Ability and Resistance to Metabolic Inhibition

720         in CD8 versus CD4 T Cells. *PLOS ONE* **9**, e104104 (2014).

721    31. Antoshechkin, A. G. On intracellular formation of ethanol and its possible role in energy

722         metabolism. *Alcohol Alcohol. Oxf. Oxfs.* **36**, 608 (2001).

723    32. Sena, L. A. *et al.* Mitochondria are required for antigen-specific T cell activation through

724        reactive oxygen species signaling. *Immunity* **38**, 225–236 (2013).

725    33. Almeida, L., Lochner, M., Berod, L. & Sparwasser, T. Metabolic pathways in T cell activation

726        and lineage differentiation. *Semin. Immunol.* **28**, 514–524 (2016).

727    34. Wang, R. & Green, D. R. Metabolic checkpoints in activated T cells. *Nat. Immunol.* **13**, 907–

728        915 (2012).

729    35. Chilmonczyk, M. A., Kottke, P. A., Stevens, H. Y., Guldberg, R. E. & Fedorov, A. G.

730        Dynamic mass spectrometry probe for electrospray ionization mass spectrometry monitoring

731        of bioreactors for therapeutic cell manufacturing. *Biotechnol. Bioeng.* **116**, 121–131 (2019).

732    36. Dwarshuis, N. J., Song, H. W., Patel, A., Kotanchek, T. & Roy, K. *Functionalized*

733        *microcarriers improve T cell manufacturing by facilitating migratory memory T cell*

734        *production and increasing CD4/CD8 ratio.* http://biorxiv.org/lookup/doi/10.1101/646760

735        (2019) doi:10.1101/646760.

736    37. NMRPipe: A multidimensional spectral processing system based on UNIX pipes |

737        SpringerLink. https://link.springer.com/article/10.1007/BF00197809.

738    38. Dieterle, F., Ross, A., Schlotterbeck, G. & Senn, H. Probabilistic quotient normalization as

739        robust method to account for dilution of complex biological mixtures. Application in 1H NMR

740        metabonomics. *Anal. Chem.* **78**, 4281–4290 (2006).

741    39. Dashti, H. *et al.* Spin System Modeling of Nuclear Magnetic Resonance Spectra for

742        Applications in Metabolomics and Small Molecule Screening. *Anal. Chem.* **89**, 12201–

743        12208 (2017).

744    40. Dashti, H. *et al.* Applications of Parametrized NMR Spin Systems of Small Molecules. *Anal.*

745        *Chem.* **90**, 10646–10649 (2018).

746    41. Holmes, E., Cloarec, O. & Nicholson, J. K. Probing Latent Biomarker Signatures and in Vivo

747        Pathway Activity in Experimental Disease States via Statistical Total Correlation

748     Spectroscopy (STOCSY) of Biofluids:  Application to HgCl2 Toxicity. *J. Proteome Res.* **5**,

749     1313–1320 (2006).

750   42. Sud, M. *et al.* Metabolomics Workbench: An international repository for metabolomics data

751     and metadata, metabolite standards, protocols, tutorials and training, and analysis tools.

752     *Nucleic Acids Res.* **44**, D463–D470 (2016).


753