1    ***N*etworks for *TR*anscriptional *A*ctivity *CE*ll a*R*rays (NTRACER)**

2    NTRACER aims to identify the dynamics of signaling processes that control an observed

3    phenotype using dynamic measurements of TFr activity[1]. NTRACER uses a combination of prior

4    knowledge and an ensemble of inference methods to determine the possible relationships

5    between the given cellular inputs and TFrs. NTRACER employs normalized activity data from

6    the significant TFrs as input. The computational pipeline involves three main steps: i) statistical

7    analysis to identify significant changes in the TFr activity data, ii) generation of an initial network

8    topology, and iii) network identification. Overall, the envisioned computational pipeline was

9    developed to identify highly robust and consistent results in the final networks by protecting

10   against the erroneous identification of edges that could result from noisy data. Robustness and

11   consistency were accomplished by a combination of data pre-filtering, structure optimization and

12   bootstrapping techniques.

13   The dynamic network was originally modeled as a three-level Boolean paradigm[1]. Herein,

14   we present an improved methodology to avoid data discretization. The most likely connections

15   present at each time are established later by minimizing the difference between the

16   experimental data and the simulation, where the relationships between the nodes are calculated

17   as linear or quadratic regression models.  Additionally, we have expanded the libraries of

18   available dynamic inference methods, we have automated the determination of the shortest

19   paths between each TFr and the applied extracellular stimuli, and we have allowed the inclusion

20   of self-loops, where a single TFr was allowed to act on itself from the TRACER experimental

21   data only. Solely new modifications added to NTRACER since our original publication[1] are

22   discussed here.

23

## Inference methods

Multiple inference methods were incorporated into NTRACER to establish new possible connections between extracellular cues (i.e., RGD and stiffness) and measured TFrs not previously reported in the scientific literature[1]. A combination of modified methods to account for dynamics, either linear or non-linear, was summarized into a unique inference network. Linear methods included PLSR[2], similarity index, SI[3], and linear ordinary differential equations (ODE) based on TIGRESS[4]. Non-linear methods included newer strategies, such as dynamic mutual information methods (ARACNE[5], CLR[6], MRNET[7], dynamic random forest[8]), as well as well-established dynamic methods, such as dynamic Bayesian networks[9]. 500 bootstrapping samples from the normalized TFr data were employed to determine connections using the above methods. If a connection appeared in any method in more than 65% of the bootstrapping cases or it appeared greater than 700 times across all methods, it was deemed significant. Those cut-offs were selected based on the frequency distributions of the bootstrapping results. The selected cut-offs coincided with initial frequency of the second distribution of the bimodal bootstrapping results (**Fig. S9**).

*Dynamic PLSR*: Dynamic PLSR was employed to infer connections between the stimuli or inputs and the TFrs using the *pls* package[10]. The first two interpolated time points (3 and 4.5 hrs) for each TFr, $X_{j,t<ti}$, were employed to regress them against the different conditions (i.e., stiffness or RGD concentration), $Y_i$ (Eq. 4). Connections were considered significant if their loads for their first component were greater than 0.15. The directionality of the interaction is given by the sign of the loading.

$$Y_i = \sum_{j=1}^{n} B_j X_{j,t<t_i} \quad \text{(Eq.4)}$$

Identification of TFrs that most likely affected other constructs was based on the differences between scaled activities acquired at two consecutive time points for the same TFr. TFr

48    activities were regressed with respect to activity at the previous time point of the constructs,

49    which is an approximation to the first derivative over time of the given TFr. Connections were

50    considered significant if their loads for their first component were greater than 0.3.

51
$$X_{i,t=t+1} - {}_{i,t=t} = \sum_{j=1}^{n} B_j X_{j,t=t} \text{ (Eq. 5)}$$

52    *Dynamic Similarity Index (SI)*: The SI is defined as the scalar product of the dynamic trajectories

53    of the average activities of two TFrs over time. Therefore, if the dynamic trends of two TFrs

54    were similar, the SI is close to 1, and if they were similar but in completely opposite directions

55    (anti-correlated), the SI value would be -1. A SI index close to 0 indicates that there is no

56    correlation between the dynamic trends of the two TFrs. Here, we calculated the SI of two

57    dynamic trajectories, but where one was delayed with respect to the other, so that we could

58    infer the directionality and sign of the observed correlation in the following manner:

59
$$SI = \frac{(average(X_{j,t+1})-1)(average(X_{i \in n,t})-1)}{\sqrt{\sum_{k=1}^{m}(average(X_{j,k})-1)^2} \sqrt{\sum_{k=1}^{m}(average(X_{j,k})-1)^2}} \text{ (Eq. 6)}$$

60    Similarly, we have employed the original definition to calculate the relationships between

61    extracellular conditions or stimuli and TFrs by only employing the first two interpolation times. All

62    the connections that have an abs(SI) ≥ 0.95 were considered significant and 0.9 in the case of

63    edges between stimuli and TFrs.

64    *ODE-TIGRESS*: Lasso regression with feature selection stabilization has been successfully

65    applied to infer biological connections[4]. Here, we presented a modification of the procedure,

66    ODE-TIGRESS, where an approximation of the first derivative over time for a given TFr is

67    regressed with respect to all the other TFrs and stimuli present in the system. Lasso regression

68    was performed using the *lars* package, with a regularization penalty, λ, equal to unity, aiming to

69    minimize L:

3

70
$$L = \frac{X_{i,t=t+1} - X_{i,t=t}}{t_{t+1} - t_t} - \sum_{j=1}^{n} \beta_j X_{j,t=t} + \lambda \sum_{j=1}^{n} \beta_j \text{ (Eq. 7)}$$

71   1000 samples were generated from the original data by randomly multiplying each value by a

72   factor between 0 to 1.  An interaction was deemed significant if it was present in at least 99% of

73   the iterations. Directionality and sign were granted by the regression parameters.

74   *Dynamic mutual information*: Mutual information (MI) methods were not only considered to

75   determine interactions between stimuli and TFrs, as in the original version of NTRACER, but

76   also between TFrs. The mutual information matrices (MIM) for relationships between inputs and

77   TFrs were constructed as for the dynamic PLSR case. The sign of each interaction between a

78   stimulus and a given TFr was determined by the initial slope over time for each stimulus. For

79   interactions between TFrs, MIM was merged from two matrices: one that contained all the data

80   except the last time point and another that contained all the data points except the first time

81   point. This method provided the MI between the different TFrs with directionality, representing

82   changes between immediately successive time points. The *minet* package[11] was selected to

83   assess the MIM with the Schurmann-Grassberger estimate of the entropy[12] by equal frequency

84   for discretization of the data for ARACNE, CLR and MRNET. Inference networks were created

85   from interactions between each TFr at an initial time point versus all TFrs at the following time

86   points with values greater than 0, as found using any of the above methods. Default parameters

87   were used otherwise.

88   *Dynamic Bayesian Networks*: Dynamic Bayesian networks were obtained assuming that all the

89   data were not independent, due to the short experimental frequency used, and no prior

90   knowledge was provided to BANJO[13], http://www.cs.duke.edu/~amink/software/banjo/.  No

91   parents were allowed for any of the stimuli, and all the data were discretized into three intervals

92   for each type of extracellular stimulus. Simulated annealing with random local moves was the

93   choice for the searching strategy with the default parameters and a maximum parent size of 5.

94   Banjo was run 500 times, and interactions were obtained from the top network for each run.

95   Interaction signs were given by the influence score.

96   *Dynamic random forest*: For the dynamic random forest version, concepts from GENIE3[14] were

97   incorporated, but with modifications to permit handling time-series data by the non-linear

98   random forest approach. The approximation to the first derivative over time was calculated as

99   above.  A total of 1000 random trees were created using the data for all the TFrs and treatments

100  for each time point employing the *randomForest* package[15]. The square root of the total number

101  of all the TFrs and conditions was used to select the number of random TFrs to start populating

102  the trees. The importance of a node was measured by the reduction in the residuals. Edges

103  were considered significant if they appeared in the top 10% ranked weights. Directionality was

104  guaranteed by the temporal order.

105  *Consensus inference network*: A total of 500 bootstrap samples were generated using the

106  weights described above and the inference methods listed above applied to each bootstrapping

107  sample. To combine all bootstrap samples, edges were deemed significant if there were present

108  in more than 65% of the runs for at least one inference method or if the number of the times that

109  was significant by some of the investigated inference methods exceeded the 700 counts (140%

110  of the 500 bootstrap samples). These cut-offs were selected based on the bimodal frequency

111  distribution for each method alone and all methods combined. Specifically, they were selected

112  to coincide with the start of the second distribution of the bimodal graph (**Fig. S9**)

113  **Determination of TFr networks evolution over time upon chemical and physical**

114  **alterations of the extracellular environment**

115     The initial network topology originated from an equally weighted number of prior knowledge

116  sources and inference methods. Prior knowledge and inference networks were combined into a

117  unique structure that served as a combined initial knowledge network model for the modified

118    version of CellNOptR[16] in NTRACER. The improved NTRACER (NTRACER v2.0) was

119    employed to identify the most likely connections present at each time point, penalizing network

120    complexity. First, the initial network was simplified by removing all connections that did not

121    include edges between the external stimuli (i.e., adhesion peptide concentration and gel

122    stiffness) and TFrs or between TFrs.

123        NTRACER v2.0 was adapted from the three-level Boolean to a continuous paradigm, where

124    edges represent linear and non-linear interactions between the nodes. This modification allowed

125    accommodating continuous variable levels (i.e., stiffness and RGD concentration). These

126    features were required in order to capture the cellular biphasic response upon chemical and

127    physical environmental cues. The prediction of the output from the model was obtained from a

128    regression model that accounts for the contributions of all the input nodes to a given TFr activity.

129    Initially the regression model was assumed linear. However, lack of fit to a linear model was

130    estimated with the rainbow test[17] ($p$-value≤0.1), and a quadratic term was added to model the

131    non-linear effects.

132        Assume that the following reactions are active in the random structure $i$:

133                                          A → B

134                                          B ⊣ C

135                                          A → C

136    NTRACER v2.0 fits a linear model for each of the output nodes, in this case, B and C, as a

137    function of their input nodes:

138
$$B_{t=t+1} = \alpha_1 A_{t=t} \quad \text{(Eq.8)}$$

139
$$C_{t=t+1} = \alpha_1 A_{t=t} - \alpha_2 B_{t=t} \quad \text{(Eq.9)}$$

140    Note that NTRACER v2.0 aims to predict the next temporal response of a given node, in this

141    case, B and C, based on the previous temporal values of A and B. In addition, for each of the

models and coefficients, NTRACER v2.0 determines the lack of fit to a linear model using the rainbow test [17] from the *lmtest* package[18]. If the alternative hypothesis is significant (p-value≤0.1), in other words, if the relationship is not linear, an additional squared term is added to the model. Assume that if $\alpha_2$ were not significant, then NTRACER v2.0 will fit the following model:

$$C_{t=t+1} = \alpha_1 A_{t=t} - \alpha_2 B_{t=t} + \alpha_3 B_{t=t}^2 \quad \text{(Eq.10)}$$

Another addition to NTRACER v2.0 is the manner in which TFrs are allowed to participate in self-loop edges. Here, we incorporated a penalty for self-loop edges and avoided models with only auto-regressive edges.

$$Score = \frac{1}{N}\left(\sum_{i=1}^{NC}(x_M - x_i)^2 + 0.1(N - NC)\right) + \frac{1}{NInp}\left(si \ e_{Pen}NSig + Stim_{Pen}SP^{(OrdT-1)}NStim + \right.$$

$$\left. InhM_{Pen}size_{Pen}NInhM + sl_{Pen}NSl\right) \text{(Eq. 11)}$$

Here, *N* is the total number of experimental observations; *NC* is the total number of simulations in which the model converged; $x_M$ represents the simulation results from the model; $x_i$ denotes the discretized experimental results; *NA*$_{Pen}$, *size*$_{Pen}$, *Stim*$_{Pen}$, *InhM*$_{Pen}$, and *sl*$_{Pen}$ are the penalties assigned to the size of non-converged simulation results, number of edges from TFrs, stimuli, InhM, and self-loops, respectively; *NInp, NSig, NStim, NInhM* and *NSl* are the size of the total number of edges, number of edges originated from TFrs, stimuli, InhM and self-loops respectively; *SP* is the stimuli policy increased to penalize the appearance of long-term stimuli edges, and *OrdT* indicates the order of the experimental time whose structure is being optimized.

Only TFrs with significantly different activities among treatments in at least one time point (meta-analysis false discovery rate (fdr)-corrected *p*-value ≤ 0.02) were subsequently studied. In order to reduce the computational time, a two-level factorial design with a central point was

165    conducted to determine the parameters that yielded the lowest score for the same number of

166    iterations (**Table S1**).

167

**References**

1. M. Weiss, B. P. Bernabé, S. Shin, S. Asztalos, S. Dubbury, M. Mui, A. Bellis, D. Bluver, D. Tonetti and J. Saez-Rodriguez, *Integrat Biol*, 2014.
2. S. Wold, M. Sjöström and L. Eriksson, *Chemometrics Intelligent Lab Sys*, 2001, **58**, 109-130.
3. A. Siletz, M. Schnabel, E. Kniazeva, A. J. Schumacher, S. Shin, J. S. Jeruss and L. D. Shea, *PLoS ONE*, 2013, **8**, e57180.
4. A.-C. Haury, F. Mordelet, P. Vera-Licona and J.-P. Vert, *BMC Sys Biol*, 2012, **6**, 145.
5. A. A. Margolin, I. Nemenman, K. Basso, C. Wiggins, G. Stolovitzky, R. D. Favera and A. Califano, *BMC Bioinformatics*, 2006, **7**, S7.
6. J. J. Faith, B. Hayete, J. T. Thaden, I. Mogno, J. Wierzbowski, G. Cottarel, S. Kasif, J. J. Collins and T. S. Gardner, *PLoS biology*, 2007, **5**, e8.
7. P. E. Meyer, K. Kontos, F. Lafitte and G. Bontempi, *EURASIP journal on bioinformatics and systems biology*, 2007, **2007**.
8. L. Breiman, *Machine learning*, 2001, **45**, 5-32.
9. J. Yu, V. A. Smith, P. P. Wang, A. J. Hartemink and E. D. Jarvis, *Bioinformatics*, 2004, **20**, 3594-3603.
10. V. A. Smith, J. Yu, T. V. Smulders, A. J. Hartemink and E. D. Jarvis, *Plos Comput Biol*, 2006, **2**, 1436-1449.
11. P. E. Meyer, F. Lafitte and G. Bontempi, *BMC Bioinformatics*, 2008, **9**, 461.
12. T. Schurmann and P. Grassberger, *Chaos*, 1996, **6**, 414-427.
13. J. Yu, V. A. Smith, P. P. Wang, A. J. Hartemink and E. D. Jarvis, *Bioinformatics*, 2004, **20**, 3594-3603.
14. V. A. Huynh-Thu, A. Irrthum, L. Wehenkel and P. Geurts, *PLoS One*, 2010, **5**.
15. A. Liaw and M. Wiener, *R News*, 2002, **2**, 18-22.
16. J. Saez-Rodriguez, L. G. Alexopoulos, J. Epperlein, R. Samaga, D. A. Lauffenburger, S. Klamt and P. K. Sorger, *Mol Sys Biol*, 2009, **5**.
17. J. M. Utts, *Commun Stat-Theory Meth*, 1982, **11**, 2801-2815.
18. A. Zeileis and T. Hothorn, *Diagnostic checking in regression relationships*, R News, 2002.

203 **Supplemental Figures and Tables**
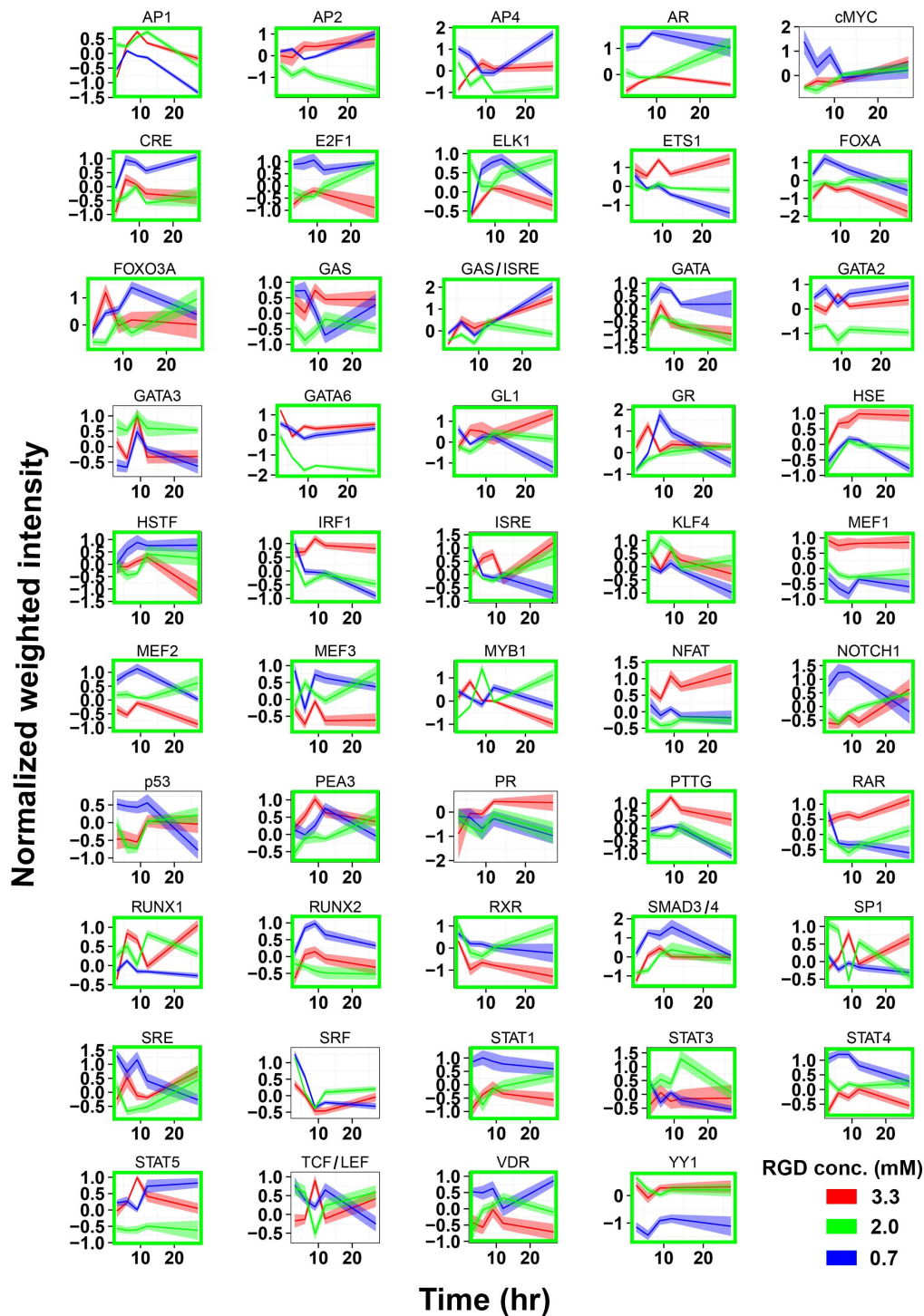


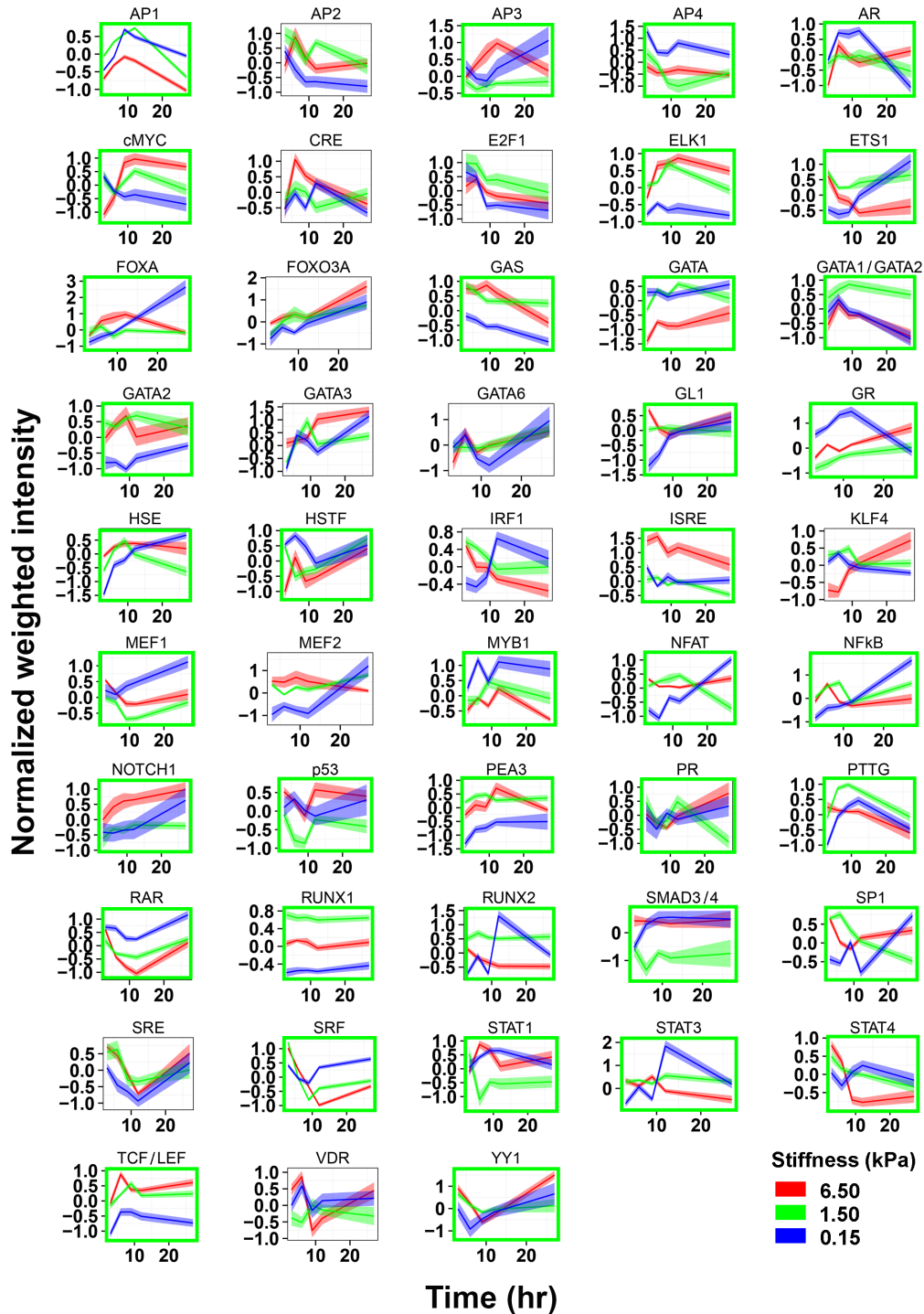| PEG (%) | 5 | 10 | 20 |
|---|---|---|---|
| I2959 (%) | 0.3 | 1 | 1 |

204

205 **Figure S1. Swelling rations of human foreskin fibroblasts cultured on PEG hydrogels**
206 **with varying modulus or RGD concentration.** Swelling ratios were calculated by comparing
207 the weight of hydrogels swollen to equilibrium (>12 hrs at room temperature) in PBS, pH 7.4
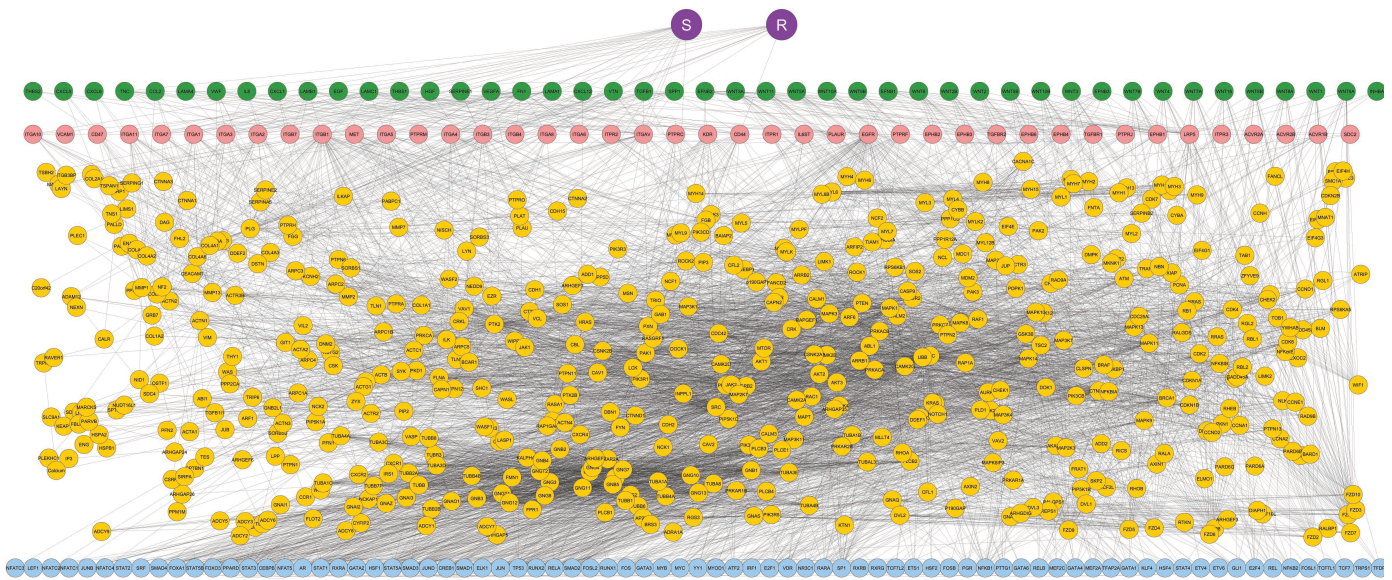208 and after lyophilization: swelling ratio = wet weight/dry weight.

209

210

**Figure S2. Dynamic TFr activity trends different levels of adhesion motif (RGD) concentration.** Mean weighted normalized and 95% confidence intervals. Green squares indicate that there is at least there is one significant difference between the RGD concentrations at one of the measured times (meta-analysis fdr-corrected p-value ≤0.02). Gel stiffness 1.5 kPa.

215

11

**Figure S3. Dynamic TFr activity trends different levels of gel stiffness.** Mean weighted
normalized and 95% confidence intervals. Green squares indicate that there is at least there is
one significant difference between the upon variation of the gel stiffness in at least one of the
measured times (meta-analysis fdr-corrected p-value ≤0.02). Adhesion concentration 2.0 mM.

12

222

**Figure S4. Mechanotransduction signaling network.** Purple nodes are the two
mechanotransduction explored variables, S for gel stiffness and R for RGD concentration; green
nodes are ligands such as fibronectin or collagen; red nodes are membrane proteins such as
receptors, integrins or cadherins; yellow nodes represent cytosolic proteins (i.e, kinases,
phosphatases); blue nodes are the transcription factors whose consensus sequences was
employed to generate the TFr employed in the study. Connections were obtained from the
GENEGO database. The initial experimental network employed for NTRACER that incorporate
the connections between the ECM and TFs also contained experimentally determined
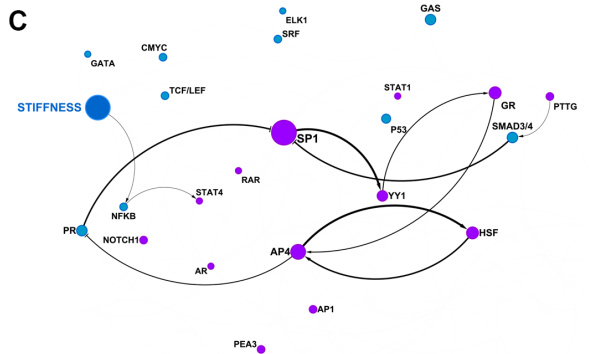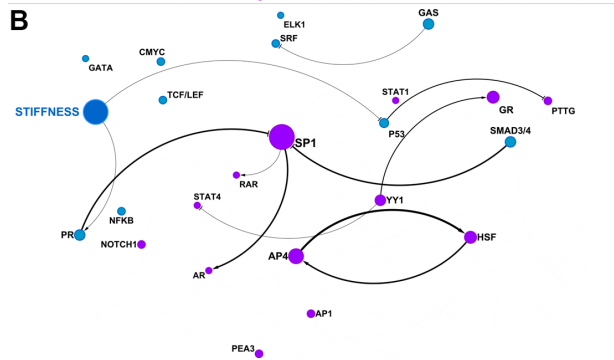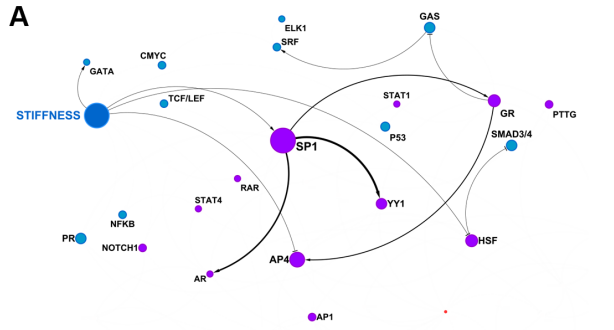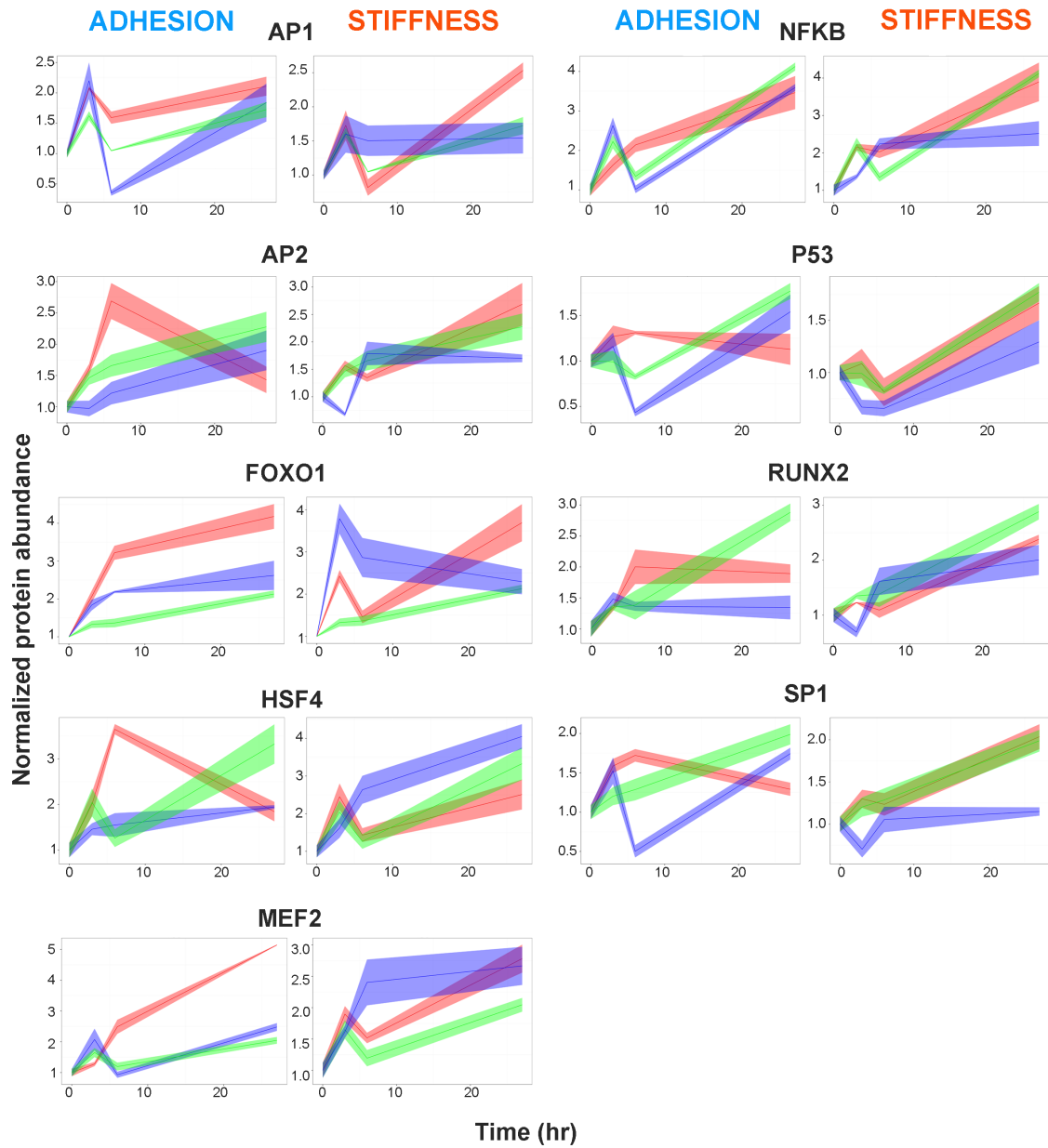connections from the adhesome database (see methods for more details).

232

233
234

235

**Figure S5 Initial networks for NTRACER for A) variations in RGD concentration; B)**
**variation in gel stiffness.** A) Number of edges originated only from literature curation (i.e., prior
knowledge), or only from inference methods and those that were common between both
approaches. Total number of edges is indicated between parenthesis. B) The two
mechanotransduction explored variables (adhesion and stiffness) correspond to the blue and
red nodes, respectively; Green nodes are transcription factors reporters. Only edges that are
common between the two sources, prior knowledge and inference methods, are represented.
Adhesion common edges are indicated in light blue and stiffness common edges are
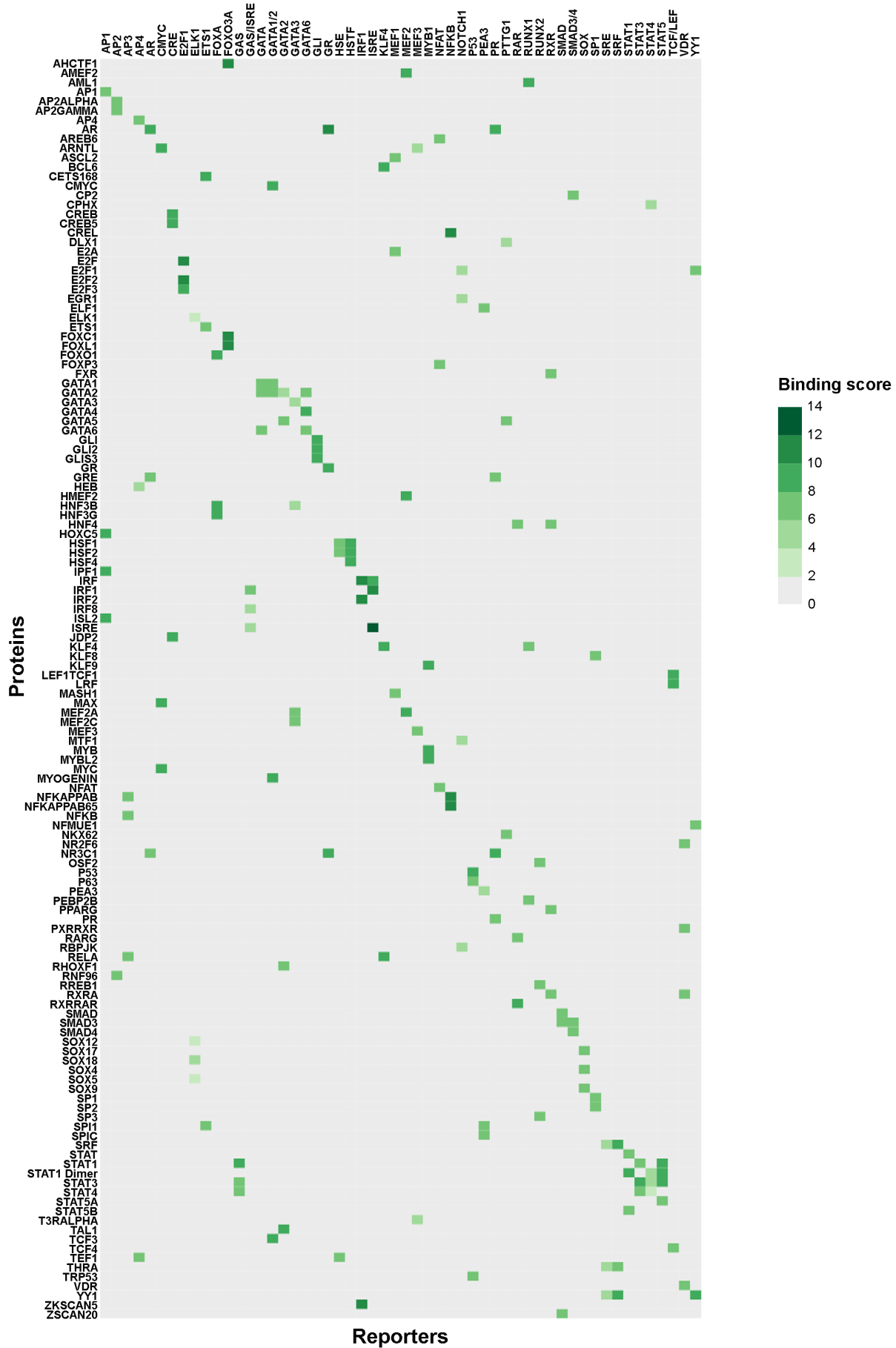represented in light red.

**Fig. S6: Dynamic TF activity networks for changes in stiffness (A-E) and RGD levels (F-J) in PEG hydrogels.** Hydrogel conditions and TFrs are represented as nodes, while the connections between them are represented by directed edges. Only edges active at each temporal step (e.g., 0-3 hrs, A and F; 3-6 hrs, B and G; 6-9 hrs, C and H; 9-12 hrs, D and I; and 12-27 hrs, E and J) are represented. Nodes affected by changes in both RGD and stiffness levels are represented in purple. Nodes only affected by RGD changes or only by stiffness changes are colored in red and aqua, respectively. Edges corresponding to linear relationships between nodes are represented by continuous lines. Edges corresponding to non-linear relationships are represented with dashes. Node size is proportional to the number of nodes that can potentially alter the TFr activity level. Similarly, edge thickness is proportional to the number of times that are activated during the measured experimental times. Activation or inhibitory effects on the downstream nodes is represented by an arrow or a T respectively.
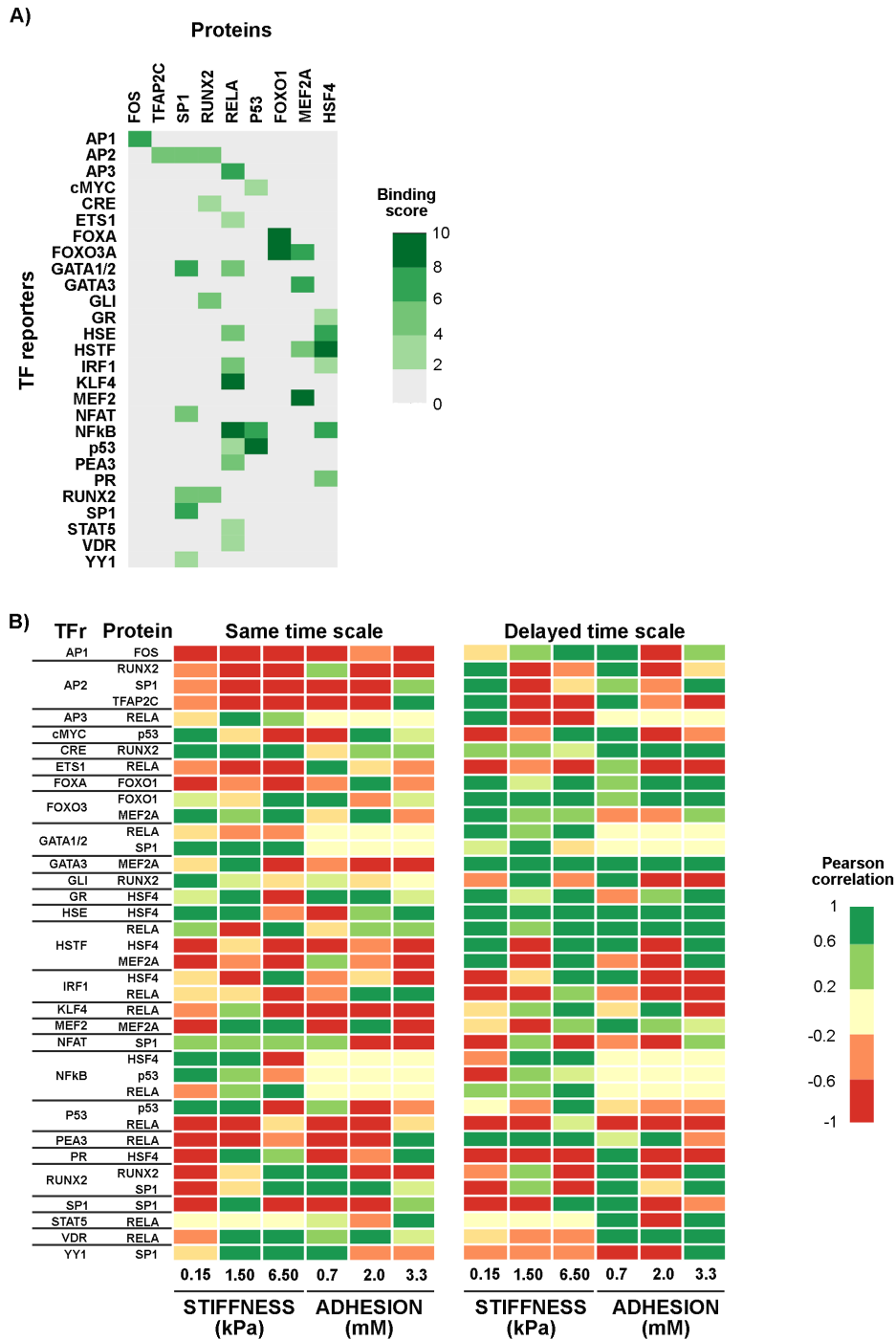
259

**Fig S7. Dynamic TF abundance trends for different levels of adhesion motif (RGD) concentration and PEG gel stiffness levels.** Mean normalized protein abundance and 95% confidence intervals. The colors of each trends are the same as Fig. S2 and Fig. S3
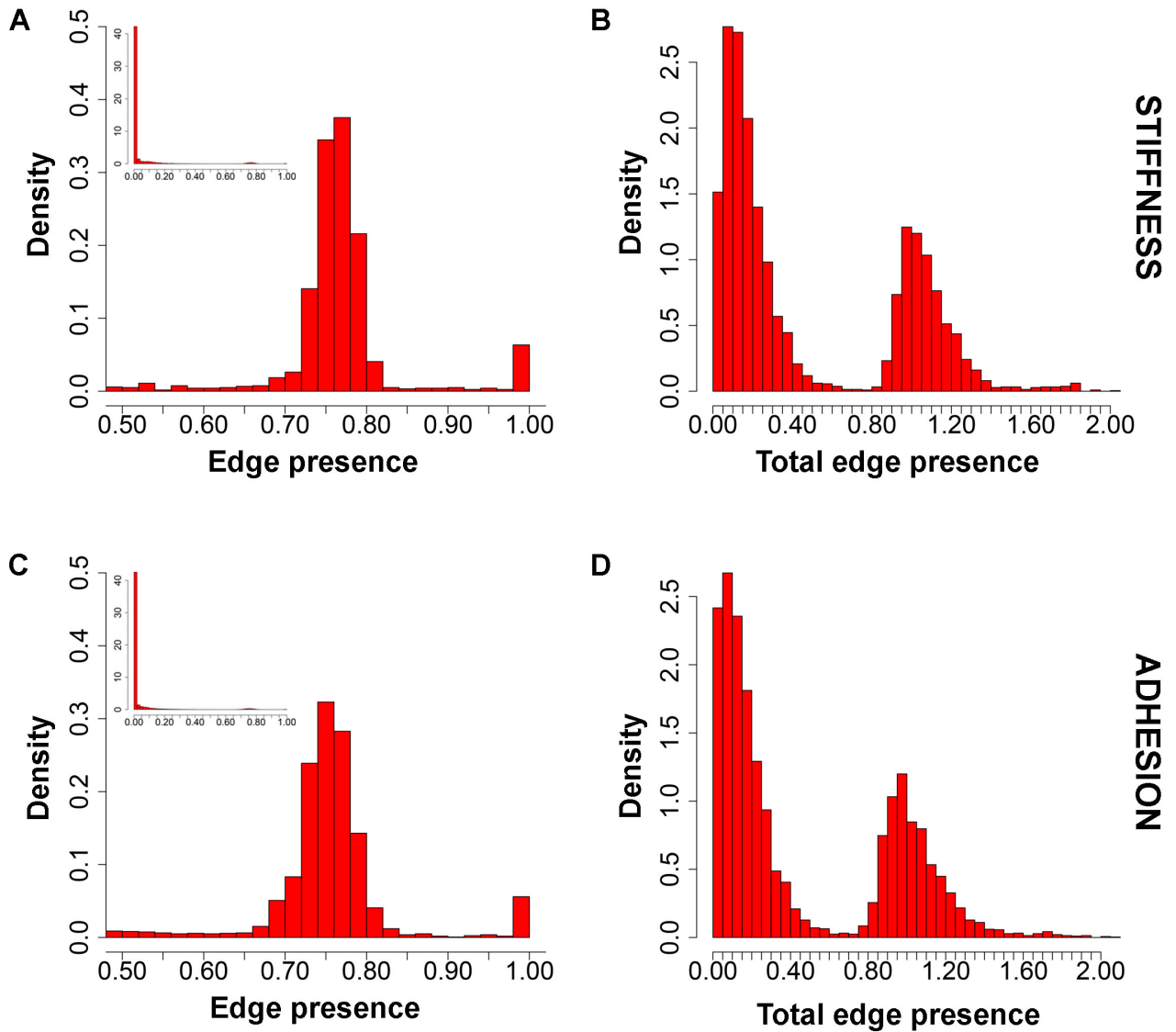
263

264

266 **Figure S8.Possible TFs binding to each of the studied TF reporters.** Each sequences for all
267 the TFr employed in the study were scanned to determine the most likely TFs whose consensus
268 binding sequences were highly similar to TFr sequence itself. The binding score represents the
269 likelihood of a given TF to bind a given reporter, accounting for sequence similarity and the non-
270 overlapping motifs. Only the top 3 rank TF for each TFr are represented as well as the TF
271 whose consensus binding sequence was employed for the design of each reporter (i.e., for AP1
272 reporter, the AP1 consensus sequence was employed).

273

**Figure S9. Validation of TRACER measurements using microWesterns arrays (MWA).** A) Possible binding sites of the proteins whose abundance was measured by MWA. Using FIMO, we identified the most likely reporters that could bind to each of the studied TFrs. We limited the list to the top 3 proteins as well as the protein that was employed for the design of the reporter (i.e., for AP1 reporter, the AP1 consensus sequence was employed). B) Most likely TF reporters that selected proteins that were analyzed by MWA arrays could bind, using the same time scale (left panel) and delayed time scale (right panel).

**Figure S10. Histograms of edge presence for the different explored inference methods.** A) Histogram for the presence of an edge in a given method (e.g., TD-PLSR, TD-MI) in the 500 bootstrapping runs for the experiments in which stiffness was altered; B) Histogram for the total summation of the presence of each edge independently of the inference method runs for the experiments in which stiffness was altered. C and D panels represents the same histograms for the experiments in which adhesion was altered.

290 **Table S1: Optimized parameters employed for NTRACER**
291

| Parameters | Values |
|---|---|
| Population size | 50 |
| Percentage of non-present edges for random start | 0.5 |
| Elitism | 5 |
| Probability of mutations | 0.001 |
| Selective pressure | 3 |
| Deactivation mechanism factor penalty | 48 |
| Edge penalty | 2 |
| Self-loop penalty | 6 |
| Stimuli penalty | 2 |

292
293
294

295    **Table S2: List of the TF antibodies employed in the microwestern arrays**

| Antibody | Company | Catalog number |
|---|---|---|
| TFAP2C/AP2-gamma | Aviva | ARP38284_T100 |
| FOS | Santa Cruz Biotechnology | sc-52 |
| FOXO1 | Cell Signaling Technologies | 9462 |
| HSF4 | Aviva | ARP32652 |
| MEF2A | Cell Signaling Technologies | 9736 |
| P53 | Cell Signaling Technologies | 9282 |
| RELA | Aviva | P100779 |
| RUNX2 | Cell Signaling Technologies | 8486 |
| SP1 | Abcam | ab13370 |
| Lamin A+C | AbCam | ab8984 |

296

297

298

299

300

301

302

303

304
305

**Table S3: List of microarrays employed for the identification of overexpressed TF gene targets in mechanotransduction related transcriptomic measurements**

| Group | Array Express ID | Cells/Tissue | Variables | FC | p-value | fdr corrected |
|-------|------------------|--------------|-----------|-----|---------|---------------|
| Stiffness | E-GEOD-22011 | Human lung fibroblasts | Different matrix stiffness | 1.3 | 0.005 | no |
| | E-GEOD-33603 | Young patient quadriceps | Massage therapy after exercise | 1.3 | 0.01 | no |
| | E-GEOD-10125 | Human dermal fibroblast cells | 3 hours of cycle mechanical loading | 1.3 | 0.01 | yes |
| RGD | E-MEXP-1273 | Human mesenchymal stem cells from adipose tissue | Monolayer or LVG or RGD alginate | 1.3 | 0.05 | no |
| Both | E-GEOD-6432 | Human fibroblasts | Culture in petri dish or attached to a tissue engineered scaffold | 1.3 | 0.01 | yes |
| | E-GEOD-44811 | Adipose stromal cells | 2D or 3D collagen culture | 1.3 | 0.01 | yes |
| | E-GEOD-39475 | Human foreskin fibroblasts | Attached versus released 3D collagen matrix | 1.3 | 0.01 | yes |
| | E-GEOD-3003 | Human CD34+ hematopoietic cells | Suspension culture or collagen I matrix | 1.2 | 0.05 | yes |
| Fibrosis | E-GEOD-17978 | Non-culture pulmonary fibroblasts from idiopathic pulmonary fibrosis (IPF) | Patients versus normal control donors | 1.3 | 0.01 | yes |