

A molecular classification of human mesenchymal stromal cells

Florian Rohart, Elizabeth A Mason, Nicholas Matigian, Rowland Mosbergen, Othmar Korn, Tyrone Chen, Suzanne Butcher, Jatin Patel, Kerry Atkinson, Kiarash Khosrotehrani, Nicholas M Fisk, Kim-Anh K-A. Lê Cao, Christine A Wells

Mesenchymal Stromal Cells (MSC) are widely used for the study of mesenchymal tissue repair, and increasingly adopted for cell therapy, despite the lack of consensus on the identity of these cells. In part this is due to the lack of specificity of MSC markers. Distinguishing MSC from other stromal cells such as fibroblasts is particularly difficult using standard analysis of surface proteins, and there is an urgent need for improved classification approaches. Transcriptome profiling is commonly used to describe and compare different cell types, however efforts to identify specific markers of rare cellular subsets may be confounded by the small sample sizes of most studies. Consequently, it is difficult to derive reproducible, and therefore useful markers. We addressed the question of MSC classification with a large integrative analysis of many public MSC datasets. We derived a sparse classifier (The Rohart MSC test) that accurately distinguished MSC from nonMSC samples with >97% accuracy on an internal training set of 635 samples from 41 studies derived on 10 different microarray platforms. The classifier was validated on an external test set of 1291 samples from 65 studies derived on 15 different platforms, with >95% accuracy. The genes that contribute to the MSC classifier formed a protein-interaction network that included known MSC markers. Further evidence of the relevance of this new MSC panel came from the high number of Mendelian disorders associated with mutations in more than 65% of the network. These result in mesenchymal defects, particularly impacting on skeletal growth and function. The Rohart MSC test is a simple *in silico* test that accurately discriminates MSC from fibroblasts, other adult stem/progenitor cell types or differentiated stromal cells. It has been implemented in the www.stemformatics.org resource, to assist researchers wishing to benchmark their own MSC datasets or data from the public domain. The code is available from the CRAN repository and all data used to generate the MSC test is available to download via the Gene Expression Omnibus or the Stemformatics resource.

1 *A Molecular Classification of Human Mesenchymal Stromal Cells*

2 Florian Rohart^{1,2}, Elizabeth Mason¹, Nicholas Matigian^{1,2}, Rowland Mosbergen^{1,5},

3 Othmar Korn¹, Tyrone Chen^{1,5}, Suzanne Butcher^{1,5}, Jatin Patel³, Kerry Atkinson³,

4 Kiarash Khosrotehrani^{3,4}, Nicholas M Fisk^{3,4}, Kim-Anh Lê Cao² and Christine A Wells^{1,5*}

5 ¹ Australian Institute for Bioengineering and Nanotechnology, The University of

6 Queensland, Brisbane, QLD Australia 4072

7 ² The University of Queensland Diamantina Institute, Translational Research Institute,

8 Woolloongabba, Brisbane QLD Australia, 4102

9 ³ The University of Queensland Centre for Clinical Research, Herston, Brisbane,

10 Queensland, Australia, 4029

11 ⁴ Centre for Advanced Prenatal Care, Royal Brisbane & Women's Hospital, Herston,

12 Brisbane, Queensland, Australia, 4029

13 ⁵ Department of Anatomy and Neuroscience, Faculty of Medicine, The University of

14 Melbourne, Victoria, 3010

15 *Correspondence to: Christine Wells, wells.c@unimelb.edu.au, @Mincle.

16 **Introduction**

17 Adult tissues maintain the capacity to be replenished as part of the normal processes of
18 homeostasis and repair. The adult stem cell hypothesis proposes that multipotent cells
19 resident in tissues are the source of this cellular renewal, and expand in response to
20 tissue injury. MSC were first isolated from bone marrow, where these occupy an
21 important stem cell niche required for reconstitution of bone and the stromal
22 compartments of marrow, and also play a supportive role in haematopoiesis
23 (Friedenstein, Piatetzky-Shapiro & Petrakova, 1966; Pittenger et al., 1999).
24 Subsequently adult stromal progenitors have been isolated and cultured from most
25 organs including placenta, heart, adipose tissue and kidneys although the identity of
26 these cells remains controversial (reviewed by (Phinney, 2012; Bianco et al., 2013)).
27 Specifically, the question of how similar cells isolated outside the bone marrow niche
28 are, and whether these could be considered bona fide MSC, or indeed challengingly,
29 whether MSC isolated from different tissues share any phenotypic or molecular
30 characteristics at all (Bianco et al., 2013). In this light various cells described as MSC
31 (whether by name or attribution) have been reported as having quite different self-
32 renewal capacity, immunomodulatory properties or propensity to differentiate *in vivo*
33 (Reinisch et al., 2014). It has been variously argued that MSC isolated from most
34 stromal tissues are derived from perivascular progenitors (Crisan et al., 2008), or
35 recruited from the bone marrow to distal tissue sites (Lee et al., 2010), or that resident
36 stromal progenitors from different tissues must have tissue-restricted phenotypes. The
37 most stringent criterion for MSC are in-vivo, bone forming capacity, however this

38 functional standard is rarely addressed in the majority of MSC studies reported in the
39 literature to date (see for example (Reinisch et al., 2014; Sworder et al., 2015)).

40

41 Several groups have attempted to address the demand for improved molecular
42 markers, for example using global proteomics methods (Li et al., 2009), epigenetic
43 markers (de Almeida et al., 2016), transcriptome analysis of cells capable of
44 regenerating the bone marrow niche (Charbord et al., 2015), or comparison of desirable
45 properties such as capacity to form bone (Sworder et al., 2015) and indeed the studies
46 reporting global 'omic analysis of MSC number in the hundreds. Each of these studies
47 identifies a different set of potential markers, but there is little consensus among them.
48 Most human studies have been conducted on very small numbers of donors, so it is
49 difficult to dissect donor-donor heterogeneity from source heterogeneity. Nevertheless,
50 variation between MSC lines is a major contributor to differences in MSC growth and
51 differentiation capacity, and clonal variation is evident even when derived from the same
52 donor bone marrow (Samsonraj et al., 2015; Sworder et al., 2015). MSC heterogeneity
53 is further compounded by growth conditions, including the density of culture, the
54 inclusion of serum, or the substrate on which they are grown (Liu et al., 2015).

55 Consequently there is little agreement in the literature on definitive molecular or cellular
56 phenotypes of human cultured MSC, whether from bone marrow or other sources.

57

58 There is little consensus on whether MSC from differing tissue sources share common
59 functional attributes. The lack of definitive markers for human MSC is a major barrier to
60 understanding genuine similarities, or resolving differences between various cell

61 sources or subsets. Even if acknowledging that there should be functional differences
62 between MSC isolated from different tissues, or donor groups, it is not clear whether
63 there should be any over-arching commonalities that might indicate shared homeostatic
64 roles or ontogenies. The field requires improved methods for benchmarking MSC
65 cultures, including molecular methods that lack the ambiguity of current antibody-based
66 methods. Here we describe a sophisticated integrative transcriptome analysis of public
67 MSC datasets, and provide a highly accurate *in silico* tool for straightforward
68 assessment of the identity of an MSC culture.

69 **Material and Methods**

70 **Design of test and training datasets**

71 A careful screening of all the datasets collated in www.stemformatics.org (Wells et al.,
72 2012), GEO (Barrett et al., 2011) and ArrayExpress (Parkinson et al., 2011) at the time
73 of this analysis identified 120 possible MSC microarray datasets. These were evaluated
74 for the availability of the primary (unprocessed) data; unambiguous replication
75 (biological not technical); the quality control metrics of RNA quality (5'-3' probe ratios);
76 linear range (box-whisker plots of sample median, min and max absolute and
77 normalized values); unambiguous sample descriptions; and sample clustering
78 concordant with the original publication. 35/120 datasets failed these criteria and were
79 excluded from the study.

80 As the range of phenotypes employed across the remaining 85 MSC microarray studies
81 was broad (Table S2), we assigned to the training group only those MSC datasets that
82 met at least the following criteria in common: Adherence, Cell surface markers CD105+,

83 CD73+, CD45- and differentiation to at least 2 of the three MSC-definitive lineages
84 (bone, cartilage or fat), and all training datasets included substantial phenotyping above
85 these minimal criteria. These minimal common criteria were hard-coded into the
86 Stemformatics annotation pipeline, we had a dedicated annotator responsible for the
87 quality of these annotations and these were reviewed independently by two additional
88 annotators. Sixteen MSC datasets met our 'gold standard' training set criteria for
89 accompanying phenotype of MSCs, together with 27 datasets containing cells from non-
90 mesenchymal or non-stromal sources, which we refer to as non-MSCs. In total, 41
91 datasets were included in the training set, with two datasets containing both MSCs and
92 non-MSCs, with a total of 125 MSC samples and 510 non-MSC samples from 10
93 different microarray platforms (Table S3, accompanies the MSC clustering in Figure 2).
94 The remaining MSC datasets were assigned to the independent test set and were used
95 only for evaluation of accuracy of the final signature.

96 Details on the samples, datasets and references of the experiments can be found in
97 Tables S2, S3 and S5. Two large datasets – 5003 (211 non-MSCs) and 6063 (45
98 MSCs), were subsampled prior to assigning to the training set to avoid unbalanced
99 results. The samples left out were included in the test set (Table S5). It consisted of 65
100 experiments (1291 samples, 213 MSCs and 499 non-MSC) profiled across 15 different
101 platforms.

102 **Pre-processing of data**

103 All data were processed using the R programming language v2.15.3 (Venables & Smith,
104 2008; R Development Core Team, 2011). The pre-processing step involved a

105 background correction performed with `affy 1.36.1` and the `affycoretools`
106 `1.30.0`, `gcrma 2.30.0`, `limma 3.14.4`, `lumi 2.10.0`, `simpleaffy 2.34.0`
107 (Gautier et al., 2004, Du, Kibbe & Lin, 2008) packages for processing of microarray data
108 depending on the platform.

109

110 Specifically, Affy GeneChips were background corrected using code:

```
111 data.bgonly <-  
112 bg.adjust.gcrma(data,affinity.info=affinity_data,fast=FALSE)  
113 ## Extract GC-RMA bg-corrected expression values without re-  
114 running additional bg-correction  
115 data.bgexpr <- rma(data.bgonly, background=FALSE,  
116 normalize=FALSE)
```

117 where:

118 "data" is loaded raw CEL data
119 "affinity_data" is precomputed probe affinity produced by
120 "compute.affinities()"

121

122 Affymetrix Gene ST arrays were RMA background corrected using `Affymetrix`
123 `Power Tools v1.14.4.1` ("`apt_probeset_summarize`" tool). Exon probe
124 expressions were summarised to the transcript level.

125

126 Illumina chips were background corrected using code:

```
127 lumiB(data, method = c('bgAdjust.affy'))
```

128 where:

129 "data" is non-normalized BeadStudio / GenomeStudio expression data returned by

130 "lumiR()"

131

132 Agilent chips were background corrected using code:

```
133 dat <- backgroundCorrect(datraw, method="normexp",
```

```
134 normexp.method="rma")
```

```
135 datbg <- dat[ dat$genes$ControlType==0, ]
```

```
136 bgave <- avereps(datbg, ID=datbg$genes[, "ProbeName"])
```

137

138 where:

139 "datraw" is non-normalised Agilent data returned by "read.maimages()"

140

141 All data was subsequently log2 transformed and a YuGene transformation was applied

142 (Lê Cao et al., 2014). YuGene is a rescaling method using the cumulative proportion

143 that is applied per sample rather than per dataset or per series. This is highly

144 advantageous as we performed 10-fold cross-validation that would otherwise require

145 renormalization as datasets were added or removed.

146

147 In order to combine all the datasets described in Table S2, probes were mapped to

148 Ensembl gene to provide a common set of identifiers. Mapping thresholds of 98% match

149 were used to align microarray probes to Ensembl human v69 transcript model cDNA

150 and ncRNA sequences obtained from Ensembl. Transcript IDs in resulting mapping

151 were converted to Gene IDs using EnsemblBiomart v69 (Zhang et al., 2011). In the

152 case of multi-mapping (several probes mapping to the same Ensembl gene ID), the
153 probe with the highest average expression was chosen, on a per-dataset basis.

154 The combined training data set included the gene expression measurement of 41,185
155 genes mapped by at least one probe in one dataset. When a dataset had no probes
156 mapping to a particular gene, the expression values of the gene were arbitrarily set to
157 zero for all samples from that dataset. A pre-screening step was then performed to
158 discard genes that were not present in at least half of the samples.

159 ***Identification of the 16-gene signature and assignation of a test sample to the***
160 ***MSC or non-MSC class***

161 The MSC signature was identified using a novel implementation of the sparse variant of
162 Partial Least Square Discriminant Analysis (sPLS-DA) ((Barker & Rayens, 2003)
163 implemented for multiple microarray studies using the mixOmics package (Lê Cao et al.,
164 2009; Lê Cao, Boitard & Philippe, 2011). Full details of the statistical model are provided
165 in the Supplementary methods. The underlying code for the statistical test is available
166 as BootsPLS in the CRAN repository, and we have also made available the d3 code for
167 the interactive MSC graph implemented in Stemformatics via the BioJS framework at
168 <http://biojs.io/d/biojs-vis-rohart-msc-test>

169

170 ***Network analysis***

171 Twenty-six genes selected on component 1 equated to 20 proteins with a curated
172 interaction in the NetworkAnalyst protein interaction database (which draws on the PPI
173 database of the International Molecular Exchange (IMEx) consortium, accessed July
174 2015 (Orchard et al., 2012; Xia, Benner & Hancock, 2014). These seed proteins were

175 annotated to a shortest-path first-order network of 36 nodes (16 seeds) and 48 PPI
176 edges. Twenty randomised sets of equivalent size were selected from the background
177 (expressed) genes to demonstrate a lack of PPI structure by chance. Gene ontology
178 analysis was assessed using hypergeometric mean against the Jan 2015 EBI UniProt
179 GO library (Huntley et al., 2015) Disease annotations were undertaken using the OMIM
180 (Baxevanis, 2012) and MGI (Shaw, 2009) databases. Subcellular location annotations
181 were taken from UniProt (EMBL, SIB Swiss Institute of Bioinformatics & Protein
182 Information Resource (PIR), 2013).

183

184 ***Differential expression analysis:***

185 Individual MSC markers were assessed for differential analysis between MSC and non-
186 MSC groups using a standard 2-tailed t-test, with a significance threshold of 10^{-6} . For
187 exploration of MSC subsets, a linear mixed model with dataset as random effect was
188 fitted for each gene for which both the mean of bone marrow samples and other sites
189 were higher than the median of all gene expression values. This retained 16,903 genes.
190 P-values were obtained by ANOVA and corrected for multiple testing with the
191 Benjamini-Hochberg procedure (Benjamini & Hochberg, 1995).

192

193 **Results**

194 ***Common MSC markers group MSC from bone marrow and other tissues.***

195 The International Society for Cellular Therapy (Dominici et al., 2006) has collated a
196 large set of markers commonly used to immunophenotype MSC. These were used, in
197 combination with more recently identified markers from the current literature (Lv et al.,

2014), to assess whether a transcript-based approach might provide a useful molecular tool to identify MSC populations (Supplementary Table S1). In order to compare data generated on different microarray platforms, we built a PLS-DA matrix using these markers and their corresponding expression in highly verified MSC samples. The resulting scatter plot (PLS-DA, Figure 1A) demonstrated the capacity to distinguish between most MSC and non-MSC samples at a transcriptional level, and further showed that MSC isolated from different tissues do cluster together using these markers. Figure 1B shows the 16 of 32 commonly used MSC markers that were significantly differentially expressed between MSC and non-MSC groups ($P < 10^{-6}$), and these included CD73 (NTE5), CD105 (Endoglin), PDGFRB and VCAM1. The average expression of the remaining markers is provided in Supplementary Figure 1. Despite ISCT recommendations, most of the MSC publications reviewed herein used a small subset of these antibodies when phenotyping MSC, and CD73+, CD105+ and CD45- were the most consistent subset used (in combination with additional markers and phenotypic information, Supplementary Table 2). When just these three markers were used to cluster all of the samples, 85% of MSC still grouped together (12/125 misclassified, Table 1, Figure 1A), but almost 12% of non-MSC samples also clustered with this group. The overall accuracy increased to 92% when all 32 markers were used, but the rate of non-MSC misclassification remained high (7%, 35/510) and the majority of these (73.5%) were fibroblasts. It may be that these markers are less stably detected at a mRNA than protein level, however this high misclassification rate is also consistent with a large body of literature documenting the ambiguity of these markers, which are shared with stromal fibroblasts, endothelial progenitors and hematopoietic cells. The

221 variable expression of all 32 markers (Figure 1B, Supplementary Figure S1) is
222 consistent with the reported variability of marker use in the wider MSC research
223 community (reviewed by (Lv et al., 2014; Samsonraj et al., 2015)). Nevertheless, the
224 capacity of these known markers to cluster MSC from different studies gave us
225 confidence that a transcriptome approach was a useful and simplified alternate to
226 antibody-based protocols, so we next took an unbiased approach to find a set of
227 markers that could improve on the current classification paradigm. Our goal was to find
228 an *in silico* marker set that reproducibly identified *bona fide* MSC samples regardless of
229 platform or laboratory differences, and provide a molecular test that was simpler, and
230 more accurate than current methods.

231
232 ***Derivation of an improved, simple and accurate in silico MSC classifier.***
233

234 A careful review of the public databases identified 120 potential MSC transcriptome
235 studies, each comprising of a small number of donors. These were carefully curated for
236 source, phenotypic information and growth conditions (see methods for details). From
237 these efforts, a gold standard ‘training set’ was identified as meeting high confidence
238 MSC phenotype including at least the minimal common set of CD73+, CD105+, CD45-
239 and bilineage differentiation. The training set consisted of 125 MSC samples from 16
240 independently derived datasets derived predominantly from bone marrow, but also
241 included studies from other adult, neonatal and fetal stromal sources. MSC were
242 compared to 510 definitively non-MSC samples from primary human tissues and cell
243 lines, including cultured fibroblasts, haematopoietic cells and pluripotent stem cell lines
244 (Supplementary Tables S2, S3).

245 To fully integrate and interrogate these data, we derived a novel cross-study analysis
246 framework. Our approach, described in Figure 2A, included a cross-platform
247 normalisation step (Lê Cao et al., 2014), and a modified variable (gene) selection
248 methodology. The first part of the protocol identified hundreds of potential MSC
249 markers, which in combination greatly improved the classification accuracy of 97.7%
250 (Table 1). This included many of the known MSC markers. Each gene was further
251 evaluated for stability by subsampling the datasets to ensure that its inclusion was not
252 reliant on one dominant source or platform. Stability is indicated by the probability of
253 selection over 200 iterations in Figure 2B, and was the step that excluded most of the
254 commonly used MSC markers. For example, PDGFRB and VCAM1 were identified as
255 potential component 1 genes but their inclusion was highly variable (0.76 and 0.59
256 probability of selection respectively).

257

258 We reasoned that if the majority of genes discriminating between MSC and non-MSC
259 are describing a common biology and are highly correlated, then a subset of these
260 genes could be identified that would represent the entire network. Therefore we
261 iteratively assessed how the inclusion of each gene contributed to the overall accuracy
262 of the signature. This found the subset of variables that were most stable and least
263 redundant at a statistical level, and that would represent the greater network of MSC-
264 related measurements (Figure 2C). Sixteen genes were identified, collectively forming a
265 'signature', which provided a high degree of discrimination between MSC and non-MSC
266 cell types, without any loss of accuracy in accurately identifying MSC (>95% correct
267 MSC call or 4/125 misclassified MSC samples, Table 1) and with improved

268 discrimination from fibroblasts and other non-MSC cell types (1.61% false positive,
269 Table 1). We confirmed that this clustering was agnostic to technology platform or
270 manufacturer (Supplementary Figure S2).

271

272 Cells derived from bone marrow were reliably grouped together with this method (Figure
273 2D, Supplementary Figure S2E), and MSC from other tissue sources, including adipose
274 tissue, skin, lung, placenta and cord blood shared this signature. Each gene in the
275 signature made an additive contribution across 4 vectors (components), such that the
276 absolute expression of any one gene might differ from sample to sample but the
277 combination of gene expression was highly predictive. High expression of component 1
278 genes was most likely to be a positive predictor of an MSC classification (Figure 2 and
279 Supplementary Figure S3A), as indicated by the correlation of expression of each gene
280 with its component. Note that the components are linear vectors, and so a negative
281 correlation (as for component 1 genes) simply indicates the contribution of the genes to
282 clustering MSC on the positive or negative region of that component. The inclusion of
283 components 2-4 provided higher discrimination for subsets of MSC and non-MSC,
284 particularly differentiating MSC and fibroblasts derived from various tissues. These latter
285 components included stress-related genes (heat shock proteins) and early indicators of
286 lineage commitment (osteomodulin). Importantly, this multicomponent based approach,
287 in contrast to a typical differential expression analysis, allowed for a common MSC
288 phenotype that is also permissive of tissue-specific differences in the wider MSC gene
289 network.

290

291 The implementation in www.stemformatics.org assessed the MSC score across 200
292 iterative predictions, where a sample must have a 95% pass rate to be classed as an
293 MSC. The distribution of the training sample scores was used to determine high
294 confidence scores (Figure 2E). By using 200 subsamplings of the training set, 200
295 scores were recorded for each sample, which enabled us to derive an individual 95%
296 Confidence Interval (CI). A sample was assigned to the MSC class if the lower bound of
297 its 95%CI is strictly higher than 0.5169. Similarly, a non-MSC classification is given if
298 the upper bound of the 95% CI was lower than 0.4337. Samples failing to meet these
299 criteria were assigned to an 'unknown' category. Accordingly, the four misclassified
300 MSC in the training set included one adult bone marrow MSC sample (predicted 1/200
301 times as MSC), and the remaining from two fetal studies, the first consisting of 10-week
302 chorionic villi (predicted 29/200 times as MSC) and 12-week chorionic membrane
303 preparation (2/200 MSC predictions), the second from a neonatal lung aspirate (0/200
304 positive MSC predictions).

305
306 ***The MSC signature genes form a cohesive network implicated in healthy***
307 ***mesenchymal development and function.***

308 To assess possible functional relationships between MSC signature genes, we used a
309 curated set of protein-protein interactions from the BioGrid database using the genes
310 selected from component 1 that showed a high discriminating power between MSC and
311 non-MSC. These formed a network of 36 interacting proteins (Figure 3A). The higher
312 expression of these genes in MSC samples is confirmed in Fig 3B. If the statistical tool
313 had identified a random set of genes, then the network would have little connectivity and
314 there would be no relevant functional annotations. This was confirmed by random

315 subsampling from the background datasets, which failed to form any PPI network. To
316 assess whether the highly connected MSC network also shared any cohesive functional
317 annotations, we examined mutation databases for evidence of human diseases
318 associated with network members. A high proportion of the MSC network (30/43) are
319 represented in Mendelian disorders of mesenchymal development by virtue of their
320 mutation spectrum in facial or musculo-skeletal dysmorphologies in man, or evidence of
321 mesodermal defects in KO mouse models (Described in detail in Supplementary Table
322 S4). These included the paired-related homeobox-1 (*PRRX1*), a transcription factor
323 important for early embryonic skeletal and facial development, and with a *de novo*
324 mutation spectrum in the embryonic dysmorphology syndrome Agnathia-otocephaly
325 (Çelik et al., 2012). Likewise, mutations in bone morphogenetic protein 14
326 (*BMP14/GDF5*) lead to developmental abnormalities in chondrogenesis and skeletal
327 bone (Degenkolbe et al., 2013). Mutations in *DDR2* cause limb defects, including
328 spondylo-epiphyseal-metaphyseal dysplasia (Ali et al., 2010) and mice over-expressing
329 *DDR2* have increased body size and atypical body fat (Kawai et al., 2014). In humans,
330 Polymorphisms in *ABI3BP* are associated with increased risk of osteochondropathy
331 (Zhang et al., 2014), and mice lacking *Abi3bp* have profound defects in MSC
332 differentiation to bone and fat (Hodgkinson et al., 2013).

333

334 We next examined functions that had been specifically validated in MSC biology,
335 specifically, whether any members of the signature had been used to prospectively
336 isolate MSC from tissue sources. *ITGA11* was a member of the core signature that has
337 been used to prospectively enrich MSC from bone marrow with enhanced colony

338 forming capacity (Kaltz et al., 2010), and independently shown to be enriched more than
339 3 fold at protein level in bone marrow MSC compared to dermal fibroblasts or
340 perivascular cells (Holley et al., 2015). Although several of the known and commonly
341 used MSC markers were indeed captured in the large initial set of potential classifiers,
342 but rejected by our statistical method on the grounds of poor selection stability, these
343 were 'rescued' in the protein interaction network. That is, the behavior of these markers
344 was variable across laboratories and between microarray platforms, and often high
345 expressed on non-MSC cell types. Nevertheless, the interaction network demonstrated
346 some cohesive biology with these known markers. The most highly connected member
347 of the extended network was VCAM1, which was identified in the large prospective
348 marker set but with a low frequency of selection (0.6 on component 1), which eliminated
349 it from the final classifier. VCAM1, together with STRO-1, has been used for the
350 prospective isolation of human bone marrow MSC (Gronthos, 2003). VCAM1 is an
351 adhesion molecule that is induced by inflammatory stimuli to regulate leukocyte
352 adhesion to the endothelium (Dansky et al., 2001); however, in cardiac precursors its
353 expression demarcates commitment to mesenchymal rather than endothelial lineages
354 (Skelton et al., 2014).

355 Other members of our network that have been previously described in human or mouse
356 MSC biology, and used to prospectively isolate cells or have been validated at the
357 protein level include *PDGFR β* (Koide et al., 2007), *SPINT2* (Roversi et al., 2014),
358 *CCDC80* (Charbord et al., 2015), *FAP* (Bae et al., 2008), *BGN* (Holley et al., 2015),
359 and *TM4SF1* (Bae et al., 2011). *SPINT2* is a serine protease inhibitor whose activity is
360 required in bone-marrow MSC, and its loss alters hematopoietic stem cell function in

361 myelo-dysplastic disorders (Roversi et al., 2014). In mouse, CCDC80 is also necessary
362 for reconstitution of bone marrow and support of haematopoiesis (Charbord et al.,
363 2015).

364

365 The network included a high proportion of extracellular proteins (54%) with
366 demonstrated roles in the modification of extracellular matrix proteins including
367 proteoglycans, as well as regulators of growth factor and cytokine signalling. This
368 included the cell migration inducing protein (KIAA1199/ CEMIP), which is secreted in its
369 mature form. It regulates Wnt and TGF β 3 signalling by depolarising hyaluronan, and
370 may alter trafficking of cytokines and growth factors to the extracellular milieu (Yoshida
371 et al., 2014). DDR2 is a receptor tyrosine kinase that interacts directly with collagens. It
372 stabilises the transcription factor SNAIL, and has been implicated in epithelial-
373 mesenchyme transitions in epithelial cancers (Zhang et al., 2013). CCDC80 binds
374 syndecan-heparin sulphate containing proteoglycans, has been shown to inhibit
375 WNT/beta-catenin signalling and has a regulatory role in adipogenesis (Tremblay et al.,
376 2009; Walczak et al., 2014). SRPX2 is a secreted chondroitin sulfate proteoglycan
377 involved in endothelial cell migration, tissue remodelling and vascular sprouting (Royer-
378 Zemmour et al., 2008). The chaperonins HSPB5/CRYAB and HSPB6 stabilise protein
379 complexes, and may assist in delivery of growth factor complexes where these are
380 present in high concentrations. In transplantation paradigms it is likely that the
381 therapeutic benefit derived from MSC is via local immunomodulatory, anti-inflammatory,
382 and/or trophic effects during the acute phase of cell therapy. The network of genes
383 identified here as enriched in MSC suggests an over-arching role for these cells in

384 modifying the extracellular environment, functions important in development as well as
385 in homeostatic regulation of adult tissues.

386

387 ***MSC differentiation, dedifferentiation and the MSC signature***

388 The majority of public microarray datasets available to us had limited phenotypic data
389 available, so these were not used to derive our MSC signature. Nevertheless we

390 annotated each of these samples as *presumptive* MSC (213 samples) or *presumptive*
391 non-MSC (499 samples) based on their origin and use in the source publication

392 (Supplementary Table S5). Where MSC were profiled during *in vitro* lineage

393 differentiation, we assigned the samples taken at intermediate time points to an

394 'unknown' category (579 samples) prior to testing these with the signature.

395 Implementation of the Rohart Test in the www.stemformatics.org resource allowed us to
396 evaluate a wide range of different experimental paradigms. Despite the lack of

397 phenotypic information associated with these datasets, the agreement between

398 publication status and our classification was high. Five percent of the presumptive non-

399 MSC (27/499) were misclassified by the signature as MSC, and around half of these

400 (>13) were neonatal or fetal dermal fibroblasts (Supplementary Table S5. Others have

401 reported MSC fractions derived from dermal tissues (reviewed in (Vaculik et al., 2012))

402 and certainly fibroblasts from other sources were not classified as MSC. Furthermore,

403 the signature could discriminate between MSC and differentiating cultures. Figure 3C

404 demonstrates loss of the MSC score during chondrogenic differentiation with the

405 addition of TGF β (Dataset 6119 (Mrugala et al., 2009)) and this pattern was

406 recapitulated for cells differentiating to mineralising bone (Data not shown, but the

407 reader is referred to the Stemformatics resource, see:
408 https://www.stemformatics.org/workbench/rohart_msc_graph?ds_id=6206#) or to
409 adipose-like cells
410 (https://www.stemformatics.org/workbench/rohart_msc_graph?ds_id=6208#) or when
411 undergoing reprogramming of an adipose-tissue derived iPSC
412 (https://www.stemformatics.org/workbench/rohart_msc_graph?ds_id=5018).

413

414 ***Comparison of MSC and adult stem/progenitor cell types***

415 The limbal cell niche hosts both limbal epithelial and stromal progenitors (Lim et al.,
416 2012), and the stromal progenitors were also classified as MSC by our tool (Dataset
417 6450). Some MSC subsets are likely to be derived from a perivascular progenitor. In our
418 hands, primary skeletal-muscle mesoangioblasts thought to be a subset of perivascular
419 cells in skeletal and smooth muscle (Dataset 6265 (Tedesco et al., 2012)), defined as
420 alkaline-phosphatase⁺ CD146⁺ CD31/Epcam⁻ CD56/Ncam⁻ with demonstrated skeletal
421 muscle differentiation, were classified as MSC (Figure 3D). In contrast, the majority of
422 cells derived from a perivascular location (and confirmed as such with tissues imaged in
423 the source publication) were not classified as MSC (Figure 3E). On examining putative
424 markers of perivascular progenitors in these samples, we could demonstrate that the
425 majority of perivascular progenitors expressed higher levels of Nestin than the majority
426 of MSC (Figure 3F). MCAM⁺ and MCAM⁻ cells were apparent in both MSC and
427 pericyte groups, although a higher proportion of perivascular progenitor expressed
428 MCAM RNA. In contrast, PDGFRA was highly expressed in MSC but not informative in
429 perivascular cells, and PDGFRB was highly expressed in both populations. Others have

430 shown that high expression of PDGFRA is associated with highly proliferative MSC
431 colonies, suggesting that its expression is associated with expansion in culture
432 (Samsonraj et al., 2015). These data are consistent with a classification hierarchy
433 determined by mouse and human lineage studies, where multipotent adult cells are
434 quiescent in a perivascular location (Crisan et al., 2008; Acar et al., 2015). Thus
435 perivascular progenitor cells with MSC differentiation capacity are defined as Rohart
436 test negative, Nestin positive in our test, and as such are distinct from a Rohart test
437 positive MSC. Cells differentiating to osteoblast, chondrocyte, adipocyte or fibroblast
438 exit the MSC state and rapidly become negative for the Rohart MSC score. Given that
439 a proportion of Rohart test positive MSC express MCAM or Nestin, the classification tool
440 may detect a phenotypic spectrum that spans the intermediates across the perivascular-
441 MSC-fibroblast hierarchy.

442

443 ***Tissue clustering of MSC is confounded by sex and MHC-1 haplotype.***

444 The capacity to group MSC-like cells is consistent with the general assumption that
445 MSC from different tissue share some common molecular properties. Many of the
446 individual studies in this reanalysis describe tissue-specific differences in MSC
447 populations. We were not able to recapitulate any of these specific differences on the
448 integrated dataset. Nevertheless, MSC from different tissues did form subclusters
449 (Supplementary Figures S2, S3), and the majority of bone marrow MSC clustered
450 together (Figure S2E). We therefore examined more broadly the genes that were
451 significantly different between bone marrow MSC and other cell types at the whole
452 transcriptome level. This analysis confirmed the observed clustering of bone marrow

453 derived MSC, distinguished by differential expression of 425 genes (adjusted $P < 0.01$,
454 Supplementary Table S6). The genes that were most differentially expressed between
455 the different MSC sources in our combined analysis were MHC class I genes, and these
456 accounted for >40% of the top 100 differentially expressed genes in the bone-marrow
457 comparisons (Supplementary Table S6). The HLA isotypes were generally, but not
458 exclusively, expressed at lower levels in bone marrow MSC (Hierarchical Cluster,
459 Supplementary Figure S3). Estrogen and progesterone receptors, and a network of
460 associated target genes were also significantly different between tissue sources
461 (Supplementary Table S6), and this may reflect a bias in the sex of the donors from
462 which tissue was sampled; although the sex of the donors was not available for a
463 majority of samples, some tissues (such as decidual sources) will be entirely female in
464 origin. Further molecular sub-classifications of MSC will therefore require much larger
465 studies that address specific clinical or differentiation properties of the cells, and must
466 also consider ascertainment biases that may introduce confounding variables such as
467 HLA subtypes or sex.

468

469 Discussion

470 Modern molecular classification tools are needed for the characterisation of MSC *ex*
471 *vivo* and *in vivo*. Antibody based methods currently rely on a subset of cell surface
472 proteins that are widely acknowledged to lack specificity, and the reliability of these
473 assays is dependant on operator expertise. Our study set out to provide an alternate
474 test that had better discrimination power than current assays, was robust and easy to
475 generate. In doing so we developed a specific gene signature that is shared by a wide-

476 variety of MSC. The “Rohart MSC test” is an *in silico* tool that has been implemented as
477 a simple online test that will be useful in standardisation or improvement of current bulk
478 isolation methods. This classification tool is available in the Stemformatics.org platform,
479 together with all the primary data used in derivation of the signature. Details on
480 submitting proprietary data to the Rohart test are available on the stemformatics.org
481 site.

482

483 All together we curated more than 120 MSC-related gene expression datasets in the
484 www.stemformatics.org resource (Wells et al., 2012); the datasets can be queried here
485 using key word, dataset ID or author, together with an implementation of the Rohart
486 MSC test.

487

488 Our approach highlights the potential robustness of biological signatures when
489 combining data from many different sources, where experimental variables such as
490 platform or batch can be reduced (Figure S2). The methods we used for derivation of a
491 common MSC classifier could be applied to the meta-analysis of any cell subset or
492 phenotype where sufficient samples can be drawn from public expression databases.

493

494 The Rohart test provides a snap shot of the current state of play in MSC biology. As an
495 *in silico* test it reflects all of the ambiguities existing in current nomenclature and culture
496 practise. We anticipate that a computational classifier will evolve as the field of MSC
497 biology evolves, and as isolation methods improve. Indeed, the question of what is an
498 MSC, and whether these are a distinct stem cell population recruited from the bone

499 marrow, as suggested by mouse studies of fetomaternal microchimerism (Seppanen et
500 al., 2013) or from perivascularity, as suggested by immunotagging of MSC-like cells
501 from perivascular regions in human tissues (Crisan et al., 2008), or are resident
502 progenitor populations specific to each organ cannot be resolved in the current study.
503 The signature itself is dependent on the quality of the MSC used in the training set. As
504 rare adult stem/progenitor cell types were under-represented in the current test or
505 training datasets, we anticipate that functional classification of MSC subtypes will
506 improve as newer sampling methods provide the means to identify and replicate these
507 cells. To highlight this point, the signature distinguishes perivascular progenitors from
508 MSC, however resolving a perivascular progenitor signature would require substantially
509 more data on this population than is currently available in the public domain. We expect
510 that further refinements in the isolation or culture of purer MSC or more precisely
511 defined functional subsets will also result in future evolutions of this *in silico* signature.

512

513 In summary, we set out to systematically review the current state of play in MSC biology
514 using a meta-analysis of transcriptome studies, and in doing so were able robustly to
515 identify a general MSC phenotype that could distinguish MSC from other cell types. The
516 resulting signature could also identify points of transition as MSC underwent
517 differentiation or reprogramming studies. Furthermore, we demonstrated that, at least at
518 a gene expression level, our *de novo* derived signature outperformed the classification
519 accuracy of the combined set of traditional MSC cell surface markers. While a signature
520 approach such as ours is not able to resolve the ontogeny or in vivo function of MSC, it
521 does provide a tool for better benchmarking and comparison of the cells grown ex vivo,

522 and will assist with comparison of cells derived for clinical purposes. The methods that
523 we describe here, and the resulting molecular classifier represent an important step
524 towards addressing the more intractable questions of MSC identity, ontogenic
525 relationships and function.

526

528 **References**

- 529 Acar M., Kocherlakota KS., Murphy MM., Peyer JG., Oguro H., Inra CN., Jaiyeola C., Zhao Z.,
530 Luby-Phelps K., Morrison SJ. 2015. Deep imaging of bone marrow shows non-dividing
531 stem cells are mainly perisinusoidal. *Nature* 526:126–130.
- 532 Ali BR., Xu H., Akawi NA., John A., Karuvantevida NS., Langer R., Al-Gazali L., Leitinger B.
533 2010. Trafficking defects and loss of ligand binding are the underlying causes of all
534 reported DDR2 missense mutations found in SMED-SL patients. *Human molecular*
535 *genetics* 19:2239–50. DOI: 10.1093/hmg/ddq103.
- 536 Bae S., Park CW., Son HK., Ju HK., Paik D., Jeon C-J., Koh GY., Kim J., Kim H. 2008.
537 Fibroblast activation protein alpha identifies mesenchymal stromal cells from human bone
538 marrow. *British journal of haematology* 142:827–830. DOI: 10.1111/j.1365-
539 2141.2008.07241.x.
- 540 Bae S., Shim SH., Park CW., Son HK., Lee HJ., Son JY., Jeon C., Kim H. 2011. Combined
541 omics analysis identifies transmembrane 4 L6 family member 1 as a surface protein marker
542 specific to human mesenchymal stem cells. *Stem cells and development* 20:197–203. DOI:
543 10.1089/scd.2010.0127.
- 544 Barker M., Rayens W. 2003. Partial least squares for discrimination. *Journal of Chemometrics*
545 17:166–173. DOI: 10.1002/cem.785.
- 546 Barrett T., Troup DB., Wilhite SE., Ledoux P., Evangelista C., Kim IF., Tomashevsky M.,
547 Marshall KA., Phillippy KH., Sherman PM., Muerter RN., Holko M., Ayanbule O., Yefanov
548 A., Soboleva A. 2011. NCBI GEO: archive for functional genomics data sets—10 years on.
549 *Nucleic Acids Res* 39:D1005–D1010. DOI: 10.1093/nar/gkq1184.
- 550 Baxevanis AD. 2012. Searching online mendelian inheritance in man (OMIM) for information on
551 genetic loci involved in human disease. *Current Protocols in Bioinformatics*. DOI:
552 10.1002/0471250953.bi0102s37.
- 553 Benjamini Y., Hochberg Y. 1995. Controlling the false discovery rate: a practical and powerful
554 approach to multiple testing. *J R Stat Soc Ser B* 57:289–300. DOI: 10.2307/2346101.
- 555 Bianco P., Cao X., Frenette PS., Mao JJ., Robey PG., Simmons PJ., Wang CY. 2013. The
556 meaning, the sense and the significance: translating the science of mesenchymal stem
557 cells into medicine. *Nat Med* 19:35–42. DOI: 10.1038/nm.3028.
- 558 Carvalho BS., Irizarry RA. 2010. A framework for oligonucleotide microarray preprocessing.
559 *Bioinformatics* 26:2363–2367. DOI: 10.1093/bioinformatics/btq431.
- 560 Çelik T., Simsek PO., Sozen T., Ozyuncu O., Utine GE., Talim B., Yiğit Ş., Boduroglu K.,
561 Kamnasaran D. 2012. PRRX1 is mutated in an otocephalic newborn infant conceived by
562 consanguineous parents. *Clinical Genetics* 81:294–297. DOI: 10.1111/j.1399-
563 0004.2011.01730.x.
- 564 Charbord P., Pouget C., Binder H., Dumont F., Stik G., Levy P., Allain F., Marchal C., Richter J.,
565 Uzan B., Pflumio F., Letourneur F., Wirth H., Dzierzak E., Traver D., Jaffredo T., Durand C.
566 2015. A Systems Biology Approach for Defining the Molecular Framework of the
567 Hematopoietic Stem Cell Niche. *Cell Stem Cell* 15:376–391. DOI:
568 10.1016/j.stem.2014.06.005.
- 569 Crisan M., Yap S., Casteilla L., Chen C-W., Corselli M., Park TS., Andriolo G., Sun B., Zheng
570 B., Zhang L., Norotte C., Teng P-N., Traas J., Schugar R., Deasy BM., Badylak S., Bühring

- 571 H-J., Giacobino J-P., Lazzari L., Huard J., Péault B. 2008. A Perivascular Origin for
572 Mesenchymal Stem Cells in Multiple Human Organs. *Cell Stem Cell* 3:301–313. DOI:
573 <http://dx.doi.org/10.1016/j.stem.2008.07.003>.
- 574 Dansky HM., Barlow CB., Lominska C., Sikes JL., Kao C., Weinsaft J., Cybulsky MI., Smith JD.
575 2001. Adhesion of monocytes to arterial endothelium and initiation of atherosclerosis are
576 critically dependent on vascular cell adhesion molecule-1 gene dosage. *Arteriosclerosis,*
577 *thrombosis, and vascular biology* 21:1662–7. DOI: 10.1161/hq1001.096625.
- 578 de Almeida DC., Ferreira MRP., Franzen J., Weidner CI., Frobel J., Zenke M., Costa IG.,
579 Wagner W. 2016. Epigenetic Classification of Human Mesenchymal Stromal Cells. *Stem*
580 *Cell Reports* 6:168–175. DOI: 10.1016/j.stemcr.2016.01.003.
- 581 Degenkolbe E., König J., Zimmer J., Walther M., Reißner C., Nickel J., Plöger F., Raspopovic
582 J., Sharpe J., Dathe K., Hecht JT., Mundlos S., Doelken SC., Seemann P. 2013. A GDF5
583 Point Mutation Strikes Twice - Causing BDA1 and SYNS2. *PLoS Genet* 9:e1003846.
- 584 Dominici M., Le Blanc K., Mueller I., Slaper-Cortenbach I., Marini F., Krause D., Deans R.,
585 Keating A., Prockop D., Horwitz E. 2006. Minimal criteria for defining multipotent
586 mesenchymal stromal cells. The International Society for Cellular Therapy position
587 statement. *Cytotherapy* 8:315–317. DOI: 10.1080/14653240600855905.
- 588 Du P., Kibbe WA., Lin SM. 2008. lumi: a pipeline for processing Illumina microarray.
589 *Bioinformatics* 24:1547–1548. DOI: 10.1093/bioinformatics/btn224.
- 590 EMBL., SIB Swiss Institute of Bioinformatics., Protein Information Resource (PIR). 2013.
591 UniProt. In: *Nucleic acids research*. 41: D43–D47.
- 592 Friedenstein AJ., Piatetzky-Shapiro II., Petrakova K V. 1966. Osteogenesis in transplants of
593 bone marrow cells. *Journal of Embryology and Experimental Morphology* 16 :381–390.
- 594 Gautier L., Cope L., Bolstad BM., Irizarry RA. 2004. affy--analysis of Affymetrix GeneChip data
595 at the probe level. *Bioinformatics* 20:307–315. DOI: 10.1093/bioinformatics/btg405.
- 596 Gronthos S. 2003. Molecular and cellular characterisation of highly purified stromal stem cells
597 derived from human bone marrow. *Journal of Cell Science* 116:1827–1835. DOI:
598 10.1242/jcs.00369.
- 599 Hodgkinson CP., Naidoo V., Patti KG., Gomez JA., Schmeckpeper J., Zhang Z., Davis B., Pratt
600 RE., Mirotsov M., Dzau VJ. 2013. Abi3bp is a multifunctional autocrine/paracrine factor that
601 regulates mesenchymal stem cell biology. *Stem Cells* 31:1669–1682. DOI:
602 10.1002/stem.1416.
- 603 Holley RJ., Tai G., Williamson AJK., Taylor S., Cain SA., Richardson SM., Merry CLR., Whetton
604 AD., Kielty CM., Canfield AE. 2015. Comparative Quantification of the Surfaceome of
605 Human Multipotent Mesenchymal Progenitor Cells. *Stem Cell Reports*. DOI:
606 10.1016/j.stemcr.2015.01.007.
- 607 Huntley RP., Sawford T., Mutowo-Meullenet P., Shypitsyna A., Bonilla C., Martin MJ.,
608 O'Donovan C. 2015. The GOA database: Gene Ontology annotation updates for 2015.
609 *Nucleic Acids Research* 43 :D1057–D1063. DOI: 10.1093/nar/gku1113.
- 610 Kaltz N., Ringe J., Holzwarth C., Charbord P., Niemeyer M., Jacobs VR., Peschel C., Häupl T.,
611 Oostendorp RAJ. 2010. Novel markers of mesenchymal stem cells defined by genome-
612 wide gene expression analysis of stromal cells from different sources. *Experimental cell*
613 *research* 316:2609–17. DOI: 10.1016/j.yexcr.2010.06.002.

- 614 Kawai I., Matsumura H., Fujii W., Naito K., Kusakabe K., Kiso Y., Kano K. 2014. Discoidin
615 domain receptor 2 (DDR2) regulates body size and fat metabolism in mice. *Transgenic*
616 *Research* 23:165–175. DOI: 10.1007/s11248-013-9751-2.
- 617 Koide Y., Morikawa S., Mabuchi Y., Muguruma Y., Hiratsu E., Hasegawa K., Kobayashi M.,
618 Ando K., Kinjo K., Okano H., Matsuzaki Y. 2007. Two distinct stem cell lineages in murine
619 bone marrow. *Stem Cells* 25:1213–1221. DOI: 2006-0325 [pii]r10.1634/stemcells.2006-
620 0325.
- 621 Lê Cao KA., Martin P., Robert-Granié C., Besse P. 2009. Sparse canonical methods for
622 biological data integration: application to a cross-platform study. *BMC Bioinformatics* 10:34.
- 623 Lê Cao K-A., Rohart F., McHugh L., Korn O., Wells CA. 2014. YuGene: A simple approach to
624 scale gene expression data derived from different platforms for integrated analyses.
625 *Genomics* 103:239–251. DOI: <http://dx.doi.org/10.1016/j.ygeno.2014.03.001>.
- 626 Lê Cao KA., Boitard S., Philippe B. 2011. Sparse PLS discriminant analysis: biologically
627 relevant feature selection and graphical displays for multiclass problems. *BMC*
628 *Bioinformatics* 12.
- 629 Lee ESM., Bou-Gharios G., Seppanen E., Khosrotehrani K., Fisk NM. 2010. Fetal stem cell
630 microchimerism: natural-born healers or killers? *Molecular Human Reproduction* 16 :869–
631 878. DOI: 10.1093/molehr/gaq067.
- 632 Li G., Zhang XA., Wang H., Wang X., Meng CL., Chan CY., Yew DTW., Tsang KS., Li K., Tsai
633 SN., Ngai SM., Han ZC., Lin MCM., He ML., Kung HF. 2009. Comparative proteomic
634 analysis of mesenchymal stem cells derived from human bone marrow, umbilical cord, and
635 placenta: Implication in the migration. *Proteomics* 9:20–30.
- 636 Lim MN., Hussin NH., Othman A., Umapathy T., Baharuddin P., Jamal R., Zakaria Z. 2012. Ex
637 vivo expanded SSEA-4+ human limb stromal cells are multipotent and do not express
638 other embryonic stem cell markers. *Molecular vision* 18:1289–300.
- 639 Liu Y., Muñoz N., Bunnell BA., Logan TM., Ma T. 2015. Density-Dependent Metabolic
640 Heterogeneity in Human Mesenchymal Stem Cells. *STEM CELLS*:n/a–n/a. DOI:
641 10.1002/stem.2097.
- 642 Lv F-J., Tuan RS., Cheung KMC., Leung VYL. 2014. Concise review: the surface markers and
643 identity of human mesenchymal stem cells. *Stem cells (Dayton, Ohio)* 32:1408–19. DOI:
644 10.1002/stem.1681.
- 645 Mrugala D., Dossat N., Ringe J., Delorme B., Coffy A., Bony C., Charbord P., Haupl T., Daures
646 JP., Noel D., Jorgensen C. 2009. Gene expression profile of multipotent mesenchymal
647 stromal cells: Identification of pathways common to TGFbeta3/BMP2-induced
648 chondrogenesis. *Cloning Stem Cells* 11:61–76. DOI: 10.1089/clo.2008.0070.
- 649 Orchard S., Kerrien S., Abbani S., Aranda B., Bhate J., Bidwell S., Bridge A., Briganti L.,
650 Brinkman FSL., Cesareni G., Chatr-aryamontri A., Chautard E., Chen C., Dumousseau M.,
651 Goll J., Hancock REW., Hannick LI., Jurisica I., Khadake J., Lynn DJ., Mahadevan U.,
652 Perfetto L., Raghunath A., Ricard-Blum S., Roechert B., Salwinski L., Stumpflen V., Tyers
653 M., Uetz P., Xenarios I., Hermjakob H. 2012. Protein interaction data curation: the
654 International Molecular Exchange (IMEx) consortium. *Nat Meth* 9:345–350.
- 655 Parkinson H., Sarkans U., Kolesnikov N., Abeygunawardena N., Burdett T., Dylag M., Emam I.,
656 Farne A., Hastings E., Holloway E., Kurbatova N., Lukk M., Malone J., Mani R., Pilicheva
657 E., Rustici G., Sharma A., Williams E., Adamusiak T., Brandizi M., Sklyar N., Brazma A.

- 658 2011. ArrayExpress update—an archive of microarray and high-throughput sequencing-
659 based functional genomics experiments. *Nucleic Acids Res* 39:D1002–D1004. DOI:
660 10.1093/nar/gkq1040.
- 661 Phinney DG. 2012. Functional heterogeneity of mesenchymal stem cells: Implications for cell
662 therapy. *J Cell Biochem* 113:2806–2812. DOI: 10.1002/jcb.24166.
- 663 Pittenger MF., Mackay AM., Beck SC., Jaiswal RK., Douglas R., Mosca JD., Moorman MA.,
664 Simonetti DW., Craig S., Marshak DR. 1999. Multilineage potential of adult human
665 mesenchymal stem cells. *Science* 284:143–147.
- 666 R Development Core Team R. 2011. R: A Language and Environment for Statistical Computing.
667 *R Foundation for Statistical Computing* 1:409. DOI: 10.1007/978-3-540-74686-7.
- 668 Reinisch A., Etchart N., Thomas D., Hofmann NA., Fruehwirth M., Sinha S., Chan CK.,
669 Senarath-Yapa K., Seo E-Y., Wearda T., Hartwig UF., Beham-Schmid C., Trajanoski S.,
670 Lin Q., Wagner W., Dullin C., Alves F., Andreeff M., Weissman IL., Longaker MT.,
671 Schallmoser K., Majeti R., Strunk D. 2014. Epigenetic and in vivo comparison of diverse
672 MSC sources reveals an endochondral signature for human hematopoietic niche formation.
673 *Blood* 125:249–260.
- 674 Roversi FM., Lopes MR., Machado-Neto JA., Longhini ALF., Duarte A da SS., Baratti MO.,
675 Palodetto B., Corrocher FA., Pericole FV., Campos P de M., Favaro P., Traina F., Saad
676 STO. 2014. Serine protease inhibitor kunitz-type 2 is downregulated in myelodysplastic
677 syndromes and modulates cell-cell adhesion. *Stem cells and development* 23:1109–20.
678 DOI: 10.1089/scd.2013.0441.
- 679 Royer-Zemmour B., Ponsole-Lenfant M., Gara H., Roll P., Lévêque C., Massacrier A., Ferracci
680 G., Cillario J., Robaglia-Schlupp A., Vincentelli R., Cau P., Szepetowski P. 2008. Epileptic
681 and developmental disorders of the speech cortex: Ligand/ receptor interaction of wild-type
682 and mutant SRPX2 with the plasminogen activator receptor uPAR. *Human Molecular
683 Genetics* 17:3617–30. DOI: 10.1093/hmg/ddn256.
- 684 Samsonraj RM., Rai B., Sathiyathan P., Puan KJ., Röttschke O., Hui JH., Raghunath M.,
685 Stanton LW., Nurcombe V., Cool SM. 2015. Establishing Criteria for Human Mesenchymal
686 Stem Cell Potency. *STEM CELLS* 33:1878–1891. DOI: 10.1002/stem.1982.
- 687 Seppanen E., Roy E., Ellis R., Bou-Gharios G., Fisk NM., Khosrotehrani K. 2013. Distant
688 mesenchymal progenitors contribute to skin wound healing and produce collagen:
689 evidence from a murine fetal microchimerism model. *PLoS One* 8:e62662. DOI:
690 10.1371/journal.pone.0062662.
- 691 Shaw DR. 2009. Searching the Mouse Genome Informatics (MGI) resources for information on
692 mouse biology from genotype to phenotype. *Current protocols in bioinformatics / editorial
693 board, Andreas D. Baxevanis ... [et al.]* Chapter 1:Unit1.7. DOI:
694 10.1002/0471250953.bi0107s25.
- 695 Skelton RJP., Costa M., Anderson DJ., Bruveris F., Finnin BW., Koutsis K., Arasaratnam D.,
696 White AJ., Rafii A., Ng ES., Elefanty AG., Stanley EG., Pouton CW., Haynes JM., Ardehali
697 R., Davis RP., Mummery CL., Elliott DA. 2014. SIRPA, VCAM1 and CD34 identify discrete
698 lineages during early human cardiovascular development. *Stem Cell Research* 13:172–
699 179. DOI: 10.1016/j.scr.2014.04.016.
- 700 Sworder BJ., Yoshizawa S., Mishra PJ., Cherman N., Kuznetsov SA., Merlino G., Balakumaran
701 A., Robey PG. 2015. Molecular profile of clonal strains of human skeletal stem/progenitor
702 cells with different potencies. *Stem cell research* 14:297–306. DOI:

- 703 10.1016/j.scr.2015.02.005.
- 704 Tedesco FS., Gerli MFM., Perani L., Benedetti S., Ungaro F., Cassano M., Antonini S.,
705 Tagliafico E., Artusi V., Longa E., Tonlorenzi R., Ragazzi M., Calderazzi G., Hoshiya H.,
706 Cappellari O., Mora M., Schoser B., Schneiderat P., Oshimura M., Bottinelli R., Sampaolesi
707 M., Torrente Y., Broccoli V., Cossu G. 2012. Transplantation of Genetically Corrected
708 Human iPSC-Derived Progenitors in Mice with Limb-Girdle Muscular Dystrophy. *Science*
709 *Translational Medicine* 4:140ra89. DOI: 10.1126/scitranslmed.3003541.
- 710 Tremblay F., Revett T., Huard C., Zhang Y., Tobin JF., Martinez R V., Gimeno RE. 2009.
711 Bidirectional modulation of adipogenesis by the secreted protein Ccdc80/DRO1/URB.
712 *Journal of Biological Chemistry* 284:8136–8147. DOI: 10.1074/jbc.M809535200.
- 713 Vaculik C., Schuster C., Bauer W., Iram N., Pfisterer K., Kramer G., Reinisch A., Strunk D.,
714 Elbe-Burger A. 2012. Human Dermis Harbors Distinct Mesenchymal Stromal Cell Subsets.
715 *J Invest Dermatol* 132:563–574.
- 716 Venables WN., Smith DM. 2008. R Development Core Team. *An Introduction to R Notes on R A*
717 *Programming Environment for Data Analysis and Graphics R core team version 2:R: A*
718 *language and environment for statistical comp.*
- 719 Walczak EM., Kuick R., Finco I., Bohin N., Hrycaj SM., Wellik DM., Hammer GD. 2014. Wnt-
720 Signaling Inhibits Adrenal Steroidogenesis by Cell-Autonomous and Non-Cell-Autonomous
721 Mechanisms. *Molecular endocrinology (Baltimore, Md.):me*20141060. DOI:
722 10.1210/me.2014-1060.
- 723 Wells CA., Mosbergen R., Korn O., Choi J., Seidenman N., Matigian NA., Vitale AM., Shepherd
724 J. 2012. Stemformatics: Visualisation and sharing of stem cell gene expression. *Stem Cell*
725 *Research* 10:387–395. DOI: <http://dx.doi.org/10.1016/j.scr.2012.12.003>.
- 726 Xia J., Benner MJ., Hancock REW. 2014. NetworkAnalyst - integrative approaches for protein-
727 protein interaction network analysis and visual exploration. *Nucleic Acids Research* 42
728 :W167–W174. DOI: 10.1093/nar/gku443.
- 729 Yoshida H., Nagaoka A., Nakamura S., Tobiishi M., Sugiyama Y., Inoue S. 2014. N-Terminal
730 signal sequence is required for cellular trafficking and hyaluronan-depolymerization of
731 KIAA1199. *FEBS letters* 588:111–6. DOI: 10.1016/j.febslet.2013.11.017.
- 732 Zhang J., Haider S., Baran J., Cros A., Guberman JM., Hsu J., Liang Y., Yao L., Kasprzyk A.
733 2011. BioMart: a data federation framework for large collaborative projects. *Database*
734 *(Oxford)* 2011:bar038. DOI: bar038 [pii]10.1093/database/bar038.
- 735 Zhang K., Corsa CA., Ponik SM., Prior JL., Piwnica-Worms D., Eliceiri KW., Keely PJ.,
736 Longmore GD. 2013. The collagen receptor discoidin domain receptor 2 stabilizes SNAIL1
737 to facilitate breast cancer metastasis. *Nature cell biology* 15:677–87. DOI:
738 10.1038/ncb2743.
- 739 Zhang F., Guo X., Zhang Y., Wen Y., Wang W., Wang S., Yang T., Shen H., Chen X., Tian Q.,
740 Tan L., Deng H-W. 2014. Genome-wide copy number variation study and gene expression
741 analysis identify ABI3BP as a susceptibility gene for Kashin–Beck disease. *Human*
742 *Genetics* 133:793–799. DOI: 10.1007/s00439-014-1418-4.

Table 1 (on next page)

Table 1

Table 1: MSC Signature improves the classification accuracy of MSC compared to a panel of 32 commonly used MSC markers. Column 1 provides the comparison of the classification accuracy of the 635 training samples using (Column 2) the 3 markers used as the minimal immunophenotype of the MSC training samples. (Column 3) a panel of 32 commonly used immune-markers in the MSC literature; (Column 4) using the unrefined sPLS-DA output; or (Column 5) with our final signature of 16 genes. Performance of each gene group was assessed using 200 random subsamplings of the training set. The internal classification error rate was calculated from a PLS-DA with 2 components (known immune-markers), or was an output of our statistical model with genes selected in an unbiased manner (cf Figure 1A).

1 **Table 1: MSC Signature improves the classification accuracy of MSC compared to**
2 **a panel of 32 commonly used MSC markers.** Column 1 provides the comparison of
3 the classification accuracy of the 635 training samples using (Column 2) the 3 markers
4 used as the minimal immunophenotype of the MSC training samples. (Column 3) a
5 panel of 32 commonly used immune-markers in the MSC literature; (Column 4) using
6 the unrefined sPLS-DA output; or (Column 5) with our final signature of 16 genes.
7 Performance of each gene group was assessed using 200 random subsamplings of the
8 training set. The internal classification error rate was calculated from a PLS-DA with 2
9 components (known immune-markers), or was an output of our statistical model with
10 genes selected in an unbiased manner (cf Figure 1A).

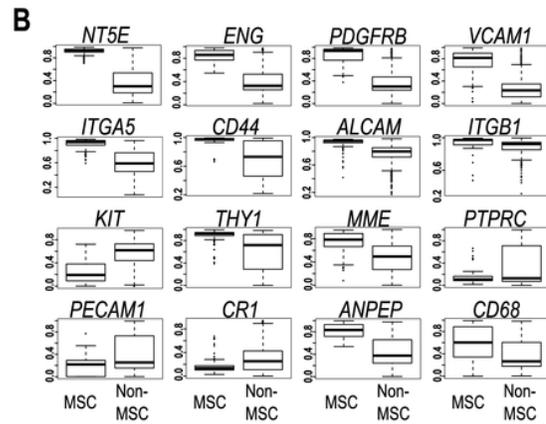
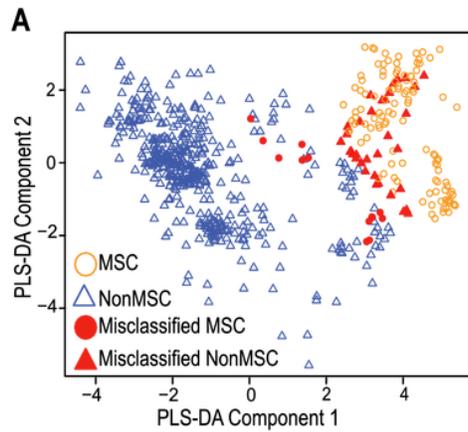
	<i>CD45,</i> <i>CD73,</i> <i>CD105</i>	32 common MSC markers	sPLS-DA prior to stable gene selection	The 16-gene MSC signature
Overall accuracy (% of 635 samples)	87.86	92.33	97.71	97.85
MSC misclassified (% of 125 samples)	14.40	11.10	3.04	4.31
Non-MSc misclassified (% of 510 samples)	11.60	6.82	2.11	1.61

11

1

Figure 1. Evaluation of Common MSC markers as transcriptional classifiers.

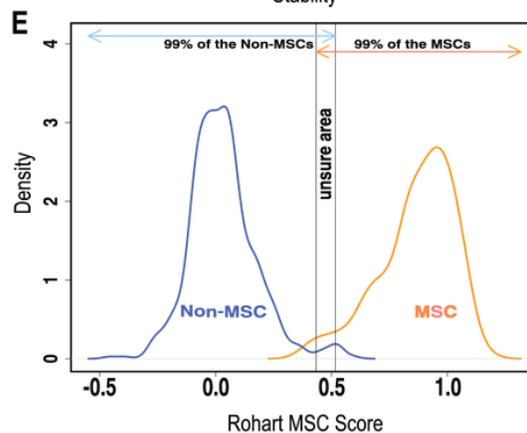
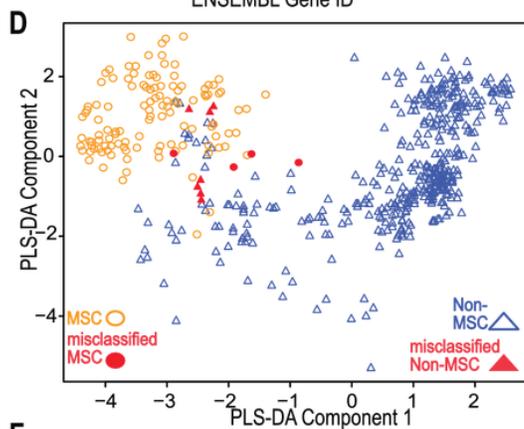
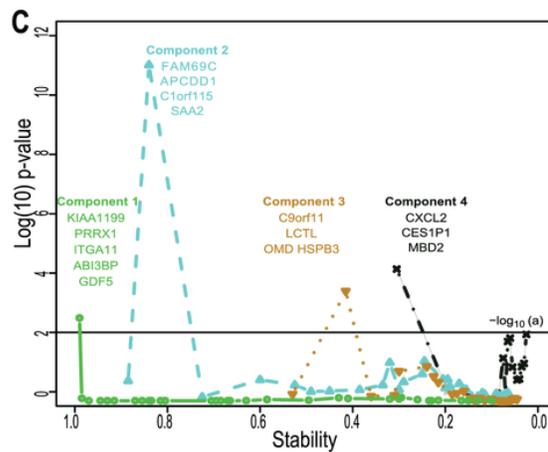
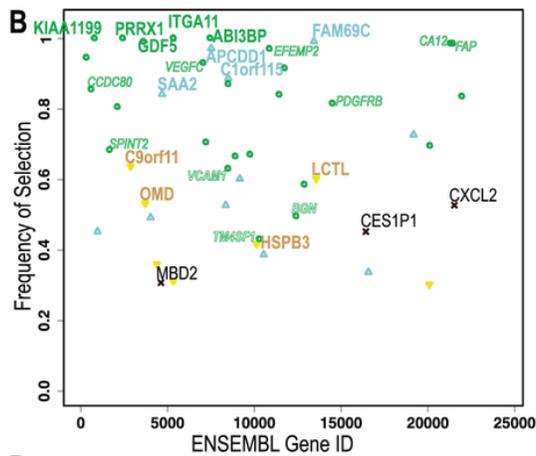
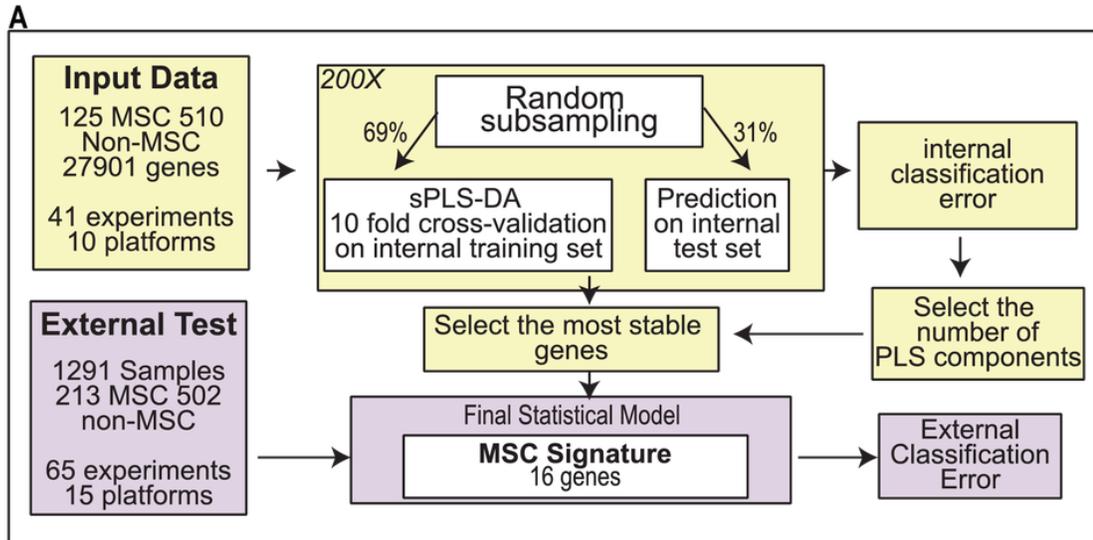
A) PLS-DA scatter plot of MSC (circles) and non-MSC cell types (triangles). Red symbols indicate cells which are incorrectly classified by the PLS-DA matrix. The matrix components consisted of 32 commonly used MSC markers. B) Box and Whisker plots showing average expression of common MSC markers that are significantly differentially expressed (t-test, $P > 10^{-6}$) between MSC (n=125) and non-MSC (n=510) cell types. See also Figure S1 and Table S1.



2

Figure 2. An improved *in silico* MSC signature.

A) Workflow summarizing the modified implementation of the sPLS-DA to integrate and evaluate cross-platform studies for derivation of a stable classifier; B) Evaluation of the stability of each gene across four components, where frequency of selection over 200 subsamplings (Y-axis) is shown per gene (ENSEMBL ID, X-axis). Labels are provided for the 16 genes contributing to the signature across 4 components. Component 1 (green), Component 2 (Blue), Component 3 (Brown), Component 4 (Black). Small text gene symbols indicate a selection of previously identified MSC markers that were excluded for poor stability. C) Evaluation of the contribution of each gene to the informativeness of its component. Each dot is a gene set, ordered along the x-axis by decreasing stability (frequency of selection). The y-axis represents the $-\log_{10}(\text{P-value})$ of a one tailed t-test indicating the improvement in classification accuracy across 4 components. D) PLS-DA scatter plot showing sample clustering and classification accuracy of the training set (635 samples) in two components (Component 1 X axis, Component 2 Y-axis). MSC samples are shown as circles, non-MSC as triangles, and misclassified samples are coloured red. E) Identifying the scores that classify an MSC or non-MSC. Distribution of the Rohart MSC Score (X-axis) and the distribution density (Y-axis) for samples in the MSC (n=115) or non-MSC (n=510) classes. Arrows indicate the scores that 99% of each class fall into. The overlap indicates the region of uncertainty, where a classification is given as 'unknown'. F) A summary of the 16-gene MSC signature colour coded to the component (as described in 1B). Gene ID is given as HUGO symbol and ENSEMBL gene ID; C is component; P is probability of selection (indicating stability); R is correlation of gene to component (as per 1D); L is predicted subcellular location of Intracellular (I), Nucleus (N), Extracellular matrix (ECM), Secreted (S), Membrane (M) and U is unknown. See also Supplementary Figure S2 and Supplemental Tables S2, S3.



F

ENSEMBL Gene ID	C	P	R	L	SYMBOL	Description
ENSG00000134046	1	1	-0.93	I, ECM	KIAA1199	Cell migration inducing protein, CEMIP
ENSG00000116132	1	1	-0.94	N	PRRX1	Paired related homeobox 1
ENSG00000137809	1	1	-0.83	M	ITGA11	Integrin alpha 11, beta 1
ENSG00000154175	1	1	-0.91	ECM	ABI3BP	ABI family member 3 (NESH) Binding Protein, TARSH
ENSG00000125965	1	0.99	-0.92	S	GDF5	Growth differentiation factor 5, BMP14
ENSG00000187773	2	0.99	-0.55	I	FAM69C	Cysteine-rich type II transmembrane protein
ENSG00000154856	2	0.96	-0.70	M	APCDD1	Adenomatosis polyposis coli down regulated 1
ENSG00000162817	2	0.89	-0.74	U (M)	C1orf115	Chromosome 1 uncharacterised open reading frame 115
ENSG00000134339	2	0.82	-0.80	S	SAA2	Serum amyloid A2
ENSG00000120160	3	0.66	0.63	I	C9orf11	Sperm acrosome associated, Equatorin
ENSG00000188501	3	0.61	-0.68	I	LCTL	Lactase-like 1
ENSG00000127083	3	0.57	-0.42	ECM	OMD	Osteomodulin
ENSG00000169271	3	0.41	0.58	N	HSPB3	Heat shock 27kDa protein 3
ENSG0000081041	4	0.53	0.54	S	CXCL2	Chemokine (C-X-C motif) ligand 2
ENSG00000134046	4	0.31	0.74	N	MBD2	Methyl-CpG binding protein 2

3

Figure 3: The MSC signature forms part of a network of extracellular proteins and discriminates between differentiating or related adult stem cell types.

A) An extended protein-protein network diagram of the Rohart MSC signature genes demonstrating a role for VCAM1 and PDGFRB as part of a functionally interconnected set of glycoproteins, integrins, growth factors and extracellular matrix proteins. Green nodes are seed network members from component 1 genes, white nodes are inferred network members, and edges are protein-protein interactions. B) Box and Whisker plots showing average expression of the genes making up the MSC signature component 1 genes in MSC (n=115) and non-MSC (n=510). Note, PRRX1, GDF5, ITGA11 and ABI3BP also form seeds in the network. KIAA1199 lacks PPI data and is not annotated in the network. C) Classification of bone marrow MSC over a time course of differentiation to cartilage; y-axis gives the Rohart score, x-axis orders the samples from each experimental series. Three differentiation series from three donors are shown. The uncertainty region stands between the MSC and non-MSC prediction regions. D) Classification of perivascular-derived stem cells from skeletal muscle mesangioblasts (HMAB), or iPSC-derived mesangioblasts (HIDEM) from donors with muscular dystrophy (MD) or healthy donors (WT). Error bars around each prediction score represent the CI boundaries. A sample is classified as 'unsure' (indicated in grey) if its prediction score or its CI overlapped the uncertainty region. E) Classification of pericytes derived from three distinct datasets: from Left-Right neonatal foreskin (Antigen HD-1 dim or bright); placental pericytes; perivascular endometrial stem cells (CD146+/PDGFRB+). Stemformatics dataset identifiers provided for each experimental series. Error bars around each prediction score represent the CI boundaries. F) Distribution of expression of common MSC/Pericyte markers. X-axis is Gene expression ranked by the YuGene cumulative proportion, Y-axis is the density distribution of MSC (orange plot, n=115) or pericytes (black plot, n= 16). See also Supplemental Figure S3 and supplemental tables S4, S5 and S6.

