

# 000 001 002 003 004 005 DENSEMARKS: LEARNING CANONICAL EMBEDDINGS 006 FOR HUMAN HEADS IMAGES VIA POINT TRACKS 007

008  
009 **Anonymous authors**  
010  
011  
012  
013  
014  
015  
016  
017  
018  
019  
020  
021  
022  
023  
024  
025  
026  
027  
028  
029  
030  
031  
032  
033  
034  
035  
036  
037  
038  
039  
040  
041  
042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053  
054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107  
108  
109  
110  
111  
112  
113  
114  
115  
116  
117  
118  
119  
120  
121  
122  
123  
124  
125  
126  
127  
128  
129  
130  
131  
132  
133  
134  
135  
136  
137  
138  
139  
140  
141  
142  
143  
144  
145  
146  
147  
148  
149  
150  
151  
152  
153  
154  
155  
156  
157  
158  
159  
160  
161  
162  
163  
164  
165  
166  
167  
168  
169  
170  
171  
172  
173  
174  
175  
176  
177  
178  
179  
180  
181  
182  
183  
184  
185  
186  
187  
188  
189  
190  
191  
192  
193  
194  
195  
196  
197  
198  
199  
200  
201  
202  
203  
204  
205  
206  
207  
208  
209  
210  
211  
212  
213  
214  
215  
216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269  
270  
271  
272  
273  
274  
275  
276  
277  
278  
279  
280  
281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300  
301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
323  
324  
325  
326  
327  
328  
329  
330  
331  
332  
333  
334  
335  
336  
337  
338  
339  
340  
341  
342  
343  
344  
345  
346  
347  
348  
349  
350  
351  
352  
353  
354  
355  
356  
357  
358  
359  
360  
361  
362  
363  
364  
365  
366  
367  
368  
369  
370  
371  
372  
373  
374  
375  
376  
377  
378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431  
432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485  
486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539  
540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
589  
590  
591  
592  
593  
594  
595  
596  
597  
598  
599  
599  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
639  
640  
641  
642  
643  
644  
645  
646  
647  
648  
649  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
699  
700  
701  
702  
703  
704  
705  
706  
707  
708  
709  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
749  
750  
751  
752  
753  
754  
755  
756  
757  
758  
759  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809  
809  
810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
859  
860  
861  
862  
863  
864  
865  
866  
867  
868  
869  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
909  
910  
911  
912  
913  
914  
915  
916  
917  
918  
919  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
969  
970  
971  
972  
973  
974  
975  
976  
977  
978  
979  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025  
1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079  
1079  
1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1129  
1130  
1131  
1132  
1133  
1134  
1135  
1136  
1137  
1138  
1139  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187  
1188  
1189  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1239  
1240  
1241  
1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1289  
1290  
1291  
1292  
1293  
1294  
1295  
1296  
1297  
1298  
1299  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349  
1349  
1350  
1351  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
1359  
1360  
1361  
1362  
1363  
1364  
1365  
1366  
1367  
1368  
1369  
1369  
1370  
1371  
1372  
1373  
1374  
1375  
1376  
1377  
1378  
1379  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387  
1388  
1389  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1399  
1400  
1401  
1402  
1403  
1404  
1405  
1406  
1407  
1408  
1409  
1409  
1410  
1411  
1412  
1413  
1414  
1415  
1416  
1417  
1418  
1419  
1419  
1420  
1421  
1422  
1423  
1424  
1425  
1426  
1427  
1428  
1429  
1429  
1430  
1431  
1432  
1433  
1434  
1435  
1436  
1437  
1438  
1439  
1439  
1440  
1441  
1442  
1443  
1444  
1445  
1446  
1447  
1448  
1449  
1449  
1450  
1451  
1452  
1453  
1454  
1455  
1456  
1457  
1458  
1459  
1459  
1460  
1461  
1462  
1463  
1464  
1465  
1466  
1467  
1468  
1469  
1469  
1470  
1471  
1472  
1473  
1474  
1475  
1476  
1477  
1478  
1479  
1479  
1480  
1481  
1482  
1483  
1484  
1485  
1486  
1487  
1488  
1489  
1489  
1490  
1491  
1492  
1493  
1494  
1495  
1496  
1497  
1498  
1499  
1499  
1500  
1501  
1502  
1503  
1504  
1505  
1506  
1507  
1508  
1509  
1509  
1510  
1511  
1512  
1513  
1514  
1515  
1516  
1517  
1518  
1519  
1519  
1520  
1521  
1522  
1523  
1524  
1525  
1526  
1527  
1528  
1529  
1529  
1530  
1531  
1532  
1533  
1534  
1535  
1536  
1537  
1538  
1539  
1539  
1540  
1541  
1542  
1543  
1544  
1545  
1546  
1547  
1548  
1549  
1549  
1550  
1551  
1552  
1553  
1554  
1555  
1556  
1557  
1558  
1559  
1559  
1560  
1561  
1562  
1563  
1564  
1565  
1566  
1567  
1568  
1569  
1569  
1570  
1571  
1572  
1573  
1574  
1575  
1576  
1577  
1578  
1579  
1579  
1580  
1581  
1582  
1583  
1584  
1585  
1586  
1587  
1588  
1589  
1589  
1590  
1591  
1592  
1593  
1594  
1595  
1596  
1597  
1598  
1599  
1599  
1600  
1601  
1602  
1603  
1604  
1605  
1606  
1607  
1608  
1609  
1609  
1610  
1611  
1612  
1613  
1614  
1615  
1616  
1617  
1618  
1619  
1619  
1620  
1621  
1622  
1623  
1624  
1625  
1626  
1627  
1628  
1629  
1629  
1630  
1631  
1632  
1633  
1634  
1635  
1636  
1637  
1638  
1639  
1639  
1640  
1641  
1642  
1643  
1644  
1645  
1646  
1647  
1648  
1649  
1649  
1650  
1651  
1652  
1653  
1654  
1655  
1656  
1657  
1658  
1659  
1659  
1660  
1661  
1662  
1663  
1664  
1665  
1666  
1667  
1668  
1669  
1669  
1670  
1671  
1672  
1673  
1674  
1675  
1676  
1677  
1678  
1679  
1679  
1680  
1681  
1682  
1683  
1684  
1685  
1686  
1687  
1688  
1689  
1689  
1690  
1691  
1692  
1693  
1694  
1695  
1696  
1697  
1698  
1699  
1699  
1700  
1701  
1702  
1703  
1704  
1705  
1706  
1707  
1708  
1709  
1709  
1710  
1711  
1712  
1713  
1714  
1715  
1716  
1717  
1718  
1719  
1719  
1720  
1721  
1722  
1723  
1724  
1725  
1726  
1727  
1728  
1729  
1729  
1730  
1731  
1732  
1733  
1734  
1735  
1736  
1737  
1738  
1739  
1739  
1740  
1741  
1742  
1743  
1744  
1745  
1746  
1747  
1748  
1749  
1749  
1750  
1751  
1752  
1753  
1754  
1755  
1756  
1757  
1758  
1759  
1759  
1760  
1761  
1762  
1763  
1764  
1765  
1766  
1767  
1768  
1769  
1769  
1770  
1771  
1772  
1773  
1774  
1775  
1776  
1777  
1778  
1779  
1779  
1780  
1781  
1782  
1783  
1784  
1785  
1786  
1787  
1788  
1789  
1789  
1790  
1791  
1792  
1793  
1794  
1795  
1796  
1797  
1798  
1799  
1799  
1800  
1801  
1802  
1803  
1804  
1805  
1806  
1807  
1808  
1809  
1809  
1810  
1811  
1812  
1813  
1814  
1815  
1816  
1817  
1818  
1819  
1819  
1820  
1821  
1822  
1823  
1824  
1825  
1826  
1827  
1828  
1829  
1829  
1830  
1831  
1832  
1833  
1834  
1835  
1836  
1837  
1838  
1839  
1839  
1840  
1841  
1842  
1843  
1844  
1845  
1846  
1847  
1848  
1849  
1849  
1850  
1851  
1852  
1853  
1854  
1855  
1856  
1857  
1858  
1859  
1859  
1860  
1861  
1862  
1863  
1864  
1865  
1866  
1867  
1868  
1869  
1869  
1870  
1871  
1872  
1873  
1874  
1875  
1876  
1877  
1878  
1879  
1879  
1880  
1881  
1882  
1883  
1884  
1885  
1886  
1887  
1888  
1889  
1889  
1890  
1891  
1892  
1893  
1894  
1895  
1896  
1897  
1898  
1899  
1899  
1900  
1901  
1902  
1903  
1904  
1905  
1906  
1907  
1908  
1909  
1909  
1910  
1911  
1912  
1913  
1914  
1915  
191

054 2019a; Cao et al., 2019) aims to follow the locations of unambiguous but isolated facial features  
 055 shared by typical faces, such as the outlines of the eyes, nose line, or mouth corners. Similarly, para-  
 056 metric 3D model estimation and tracking (Blanz & Vetter, 1999; Li et al., 2017b; Dai et al., 2020)  
 057 assumes that most face and head geometries follow a statistical shape model and can be represented  
 058 by a shared, comparatively simple mesh template.

059 However, features like hair, accessories, and clothing are often omitted from head tracking, which  
 060 typically focuses on landmarks or skin. In a typical video capture, individual landmarks or entire  
 061 head regions easily become occluded due to an extreme pose, expression, or a worn accessory,  
 062 introducing large errors in tracking. As a result, head tracks produced by conventional approaches  
 063 are fundamentally limited by their incompleteness and correspondence instability.

064 To improve robustness of correspondence search, one path forward is to extract and match repre-  
 065 sentations densely in each image pixel instead of detection and alignment of isolated landmarks.  
 066 Recent image-based vision foundational models (VFM) are one suitable source of such dense repre-  
 067 sentations known to be effective in many vision tasks (Dutt et al., 2024; Siméoni et al., 2025).  
 068 As human heads constitute a visual category with high structural similarity across instances, it is  
 069 natural to expect such representations defined at unambiguous facial features to be nearly view- and  
 070 time-invariant, facilitating exact correspondence search.

071 Building on these insights, we propose DenseMarks, a new learned representation for human heads  
 072 designed to (1) enable high-quality dense correspondences for complete human heads, including ir-  
 073 regular features such as hair or accessories, (2) achieve robust tracking under challenging conditions  
 074 such as strong occlusions, and (3) produce a structured, interpretable, and smooth canonical latent  
 075 space for exploration and interaction. We use a ViT neural backbone to predict dense per-pixel repre-  
 076 sentations within the head mask of an input image; leveraging powerful pre-trained VFM (Siméoni  
 077 et al., 2025). These representations are projected into a shared 3D space, reducing correspondence  
 078 to nearest-neighbor search and enabling intuitive interactions (e.g., click-based retrieval). To train  
 079 without ground-truth dense correspondences, we construct a diverse dataset of human head videos  
 080 with 2D point tracks from an off-the-shelf tracker (Karaev et al., 2024a). We enforce fine-grained  
 081 cross-subject consistency by optimizing a contrastive loss on matched pairs, and integrate semantic  
 082 and smoothness constraints to structure the latent space and improve interpretability.

083 We benchmark against pre-trained VFM variants (Siméoni et al., 2025; Khirodkar et al., 2024; Yue  
 084 et al., 2024), with assessment focused on dense image warping and geometric consistency measures.

## 086 2 RELATED WORK

088 **Face, Head, and Full Body Tracking.** Commonly, tracking humans in videos involves extracting  
 089 relevant information for the estimation and alignment of their pose and shape. In the simplest form,  
 090 this is achieved by predicting locations of characteristic landmark points with fixed semantics (Sag-  
 091 onas et al., 2013; Moon et al., 2020; Jin et al., 2020) using learned models (Bulat & Tzimiropoulos,  
 092 2017; Lugaressi et al., 2019a; Cao et al., 2019; Simon et al., 2017; Li et al., 2022). Ease of collecting  
 093 annotations and efficiency of landmark detectors have made landmarks essential in practical tracker  
 094 design, enabling initial rigid alignment (Qian, 2024; Qian et al., 2024; Grassal et al., 2021; Bogo  
 095 et al., 2016; Kanazawa et al., 2018; Kocabas et al., 2020). However, relying on a finite number  
 096 of isolated, sparse landmarks can compromise robustness, commonly requiring regularization or  
 097 postprocessing such as temporal smoothing (Qian, 2024; Zielonka et al., 2022; Zheng et al., 2023a;  
 098 Huang et al., 2022; Jiang et al., 2022).

099 Many methods for estimating and tracking parametric models of faces and bodies (3DMMs (Blanz &  
 100 Vetter, 1999; Zhu et al., 2017; Li et al., 2017b; Zhang et al., 2023b; Romero et al., 2017; Loper et al.,  
 101 2015; Dai et al., 2020)) are based on the *analysis-by-synthesis* paradigm (Blanz & Vetter, 1999; Zhu  
 102 et al., 2017; Feng et al., 2021; Zielonka et al., 2022; Daněček et al., 2022) that involves a combination  
 103 of rigid alignment and optimization of denser losses. While offering higher geometric completeness,  
 104 such models rely on a simple mesh topology and a limited range of geometries captured by a PCA  
 105 basis (Abdi & Williams, 2010; Jolliffe, 2011); for fitting, they commonly depend on prior landmarks  
 106 estimation and optimize highly non-convex (e.g., photometric or depth) losses.

107 Our method naturally complements 3DMM-based head trackers by supplying dense, robust semantic  
 108 correspondences for complete heads and includes features not trivially captured by landmarks or

108 parametric models (e.g., hair). This idea is similar to works that learn to predict texture coordinates  
 109 for alignment of parametric face (Feng et al., 2018; Giebenhain et al., 2025) and body (Güler et al.,  
 110 2018; Ianina et al., 2022) models, or compute multi-dimensional features, normals, and depth using  
 111 foundation models optimized for the human domain (Khirodkar et al., 2024).

112 **Canonical Space Learning.** Our method represents input samples by learned embeddings in a  
 113 shared (*canonical*) space. The idea of using canonical representations for category-level object local-  
 114 ization and pose estimation was pioneered by Normalized Object Coordinate Space (NOCS) (Wang  
 115 et al., 2019) and subsequently extended to handle sparse views, lack of dense labels, or multiple cat-  
 116 egories (Min et al., 2023; Xu et al., 2024; Krishnan et al., 2024). However, directly learning NOCS  
 117 representations for 3D heads is difficult as large collections of 3D models are absent in the human  
 118 head domain.

119 Shape correspondence task can be formulated as a problem of finding a mapping between spaces  
 120 of functions defined on shapes (Ovsjanikov et al., 2012; Rodolà et al., 2017). Existing methods  
 121 applying such functional maps for finding full-body correspondences (Neverova et al., 2020; Ianina  
 122 et al., 2022) require fitting parametric 3D models for supervision. To enable modeling parts of  
 123 human heads absent from parametric models, we opted not to use these in our training.

124 The idea of using canonical space is widespread in 3D-aware per-scene human fitting (Gafni et al.,  
 125 2021; Park et al., 2021) and human generative modeling EG3D (Chan et al., 2022; Dong et al.,  
 126 2023). Similarly, several works focus on producing unsupervised shape correspondences, in part  
 127 based on functional maps (Halimi et al., 2019; Cao & Bernard, 2022; Cao et al., 2023; Liu et al.,  
 128 2025).

129 **Embeddings from Foundation Models.** Recent progress in ViT-based VFM (Caron et al., 2021;  
 130 Oquab et al., 2023; Siméoni et al., 2025; Weinzaepfel et al., 2022; Dosovitskiy et al., 2020; Han  
 131 et al., 2022) and evidence of their emerging understanding of 3D world (Zhang et al., 2024b; Sucar  
 132 et al., 2025; Chen et al., 2025a) has fueled efforts to improve their 3D-awareness through fine-  
 133 tuning (Yue et al., 2024; Zhang et al., 2024a). Similarly, directly training siamese ViT networks on  
 134 pairs of stereo views has been shown to efficiently establish dense correspondences (Wang et al.,  
 135 2024; Leroy et al., 2024; Smart et al., 2024; Chen et al., 2025b), when prompted with 2+ images.

136 Another class of VFM, pre-trained diffusion models (e.g., Stable Diffusion (Rombach et al., 2021)),  
 137 allow inferring semantic correspondences from their image-based representations (Hedlin et al.,  
 138 2023; Zhang et al., 2023a; Zhu et al., 2024) that could be distilled into dense surface correspon-  
 139 dences across objects of arbitrary categories (Dutt et al., 2024). In our experiments, we found the  
 140 correspondences arising from point tracking (cf. next paragraph) more reliable than those arising  
 141 from pretrained diffusion models. Our method benefits from integrating VFM as a feature extractor;  
 142 in contrast to generic pre-trained deep features correlated with visual semantics, our geometry-aware  
 143 representations yield an interpretable 3D canonical space.

144 **Point Tracking.** The advent of talking heads datasets (Wang et al., 2021; Zhu et al., 2022; Ephrat  
 145 et al., 2018) and point trackers calls for approaches to tracking faces and bodies, free of an underly-  
 146 ing coarse parametric model. In particular, in a line of works starting from PIPs (Harley et al., 2022),  
 147 deep learning based methods are proposed to track any queried point along the video. Progress in  
 148 the area of point trackers has been additionally accelerated by the appearance of suitable bench-  
 149 marks, such as Tap-Vid (Doersch et al., 2022) and PointOdyssey (Zheng et al., 2023b). A series  
 150 of consequent improvements of track-any-point algorithms (Doersch et al., 2023; Li et al., 2024;  
 151 Cho et al., 2024) led to the emerging branch of CoTracker works (Karaev et al., 2024b;a), as well  
 152 as BootsTAP (Doersch et al., 2024). Similarly, a few methods rely on foundation models, such as  
 153 DINO-tracker (Tumanyan et al., 2024) for tracking any point or VGGT (Wang et al., 2025) that uses  
 154 point tracks for 3D understanding. Applications of modern algorithmic ideas for point tracking also  
 155 led to the appearance of simultaneous reconstruction and tracking methods such as Dynamic 3D  
 156 Gaussians (Luiten et al., 2024), St4rTrack (Feng et al., 2025), or Tracks-to-4D (Kasten et al., 2024).  
 157 For the downstream tasks of human tracking, similar to our method, some of the recent approaches  
 158 also make use of point tracking (Kim et al., 2025; Taubner et al., 2024) or motion data (Shin et al.,  
 159 2024).

160  
 161

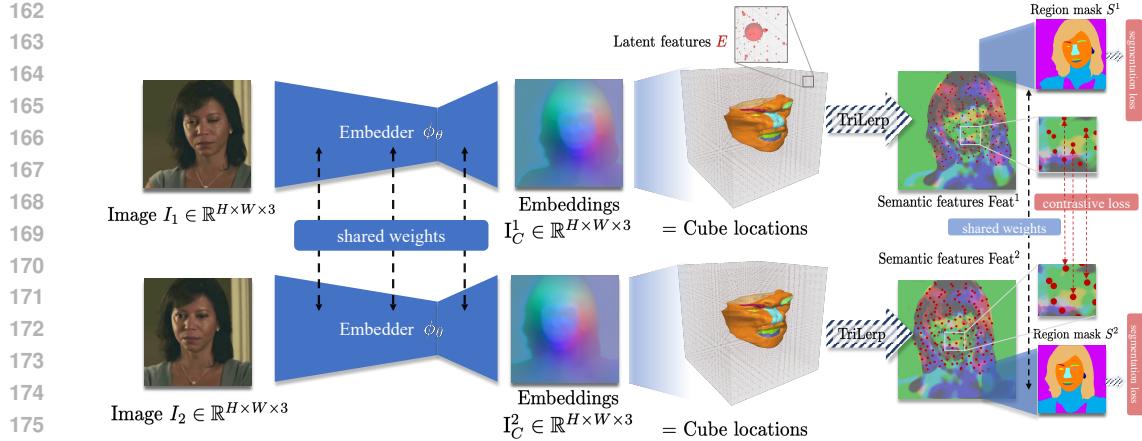


Figure 2: To learn our representation, we train an embedder network  $\phi_\theta$  in a siamese fashion. By feeding two image frames from a talking head video of the same person into the embedder independently, we obtain DenseMarks embeddings  $I_C^1, I_C^2$ . These embeddings correspond to canonical locations in the unit cube (DenseMarks space). This cube is discretized in advance, and a learnable matrix  $E$  of latent features represents  $D$ -dimensional vectors, storing semantic info of each of the voxel grid locations. To transform each of the estimated cube locations into semantic features  $\text{Feat}^1, \text{Feat}^2$ , we query  $E$  at locations  $I_C^1, I_C^2$  via trilinear interpolation (TriLerp). For the images  $I_1, I_2$ , we have a set of pair matches  $K_{\text{gt}}^1, K_{\text{gt}}^2$ , estimated by an off-the-shelf point tracker (Karaev et al., 2024a). We apply contrastive loss (Radford et al., 2021) to the semantic features of images in these locations. This way, the cube locations corresponding to the same semantic feature are pushed closer together. Additionally, we estimate region masks  $S^1, S^2$  by a semantic network  $S_\xi$  and apply segmentation loss.

### 3 METHOD

In this section, we define the representation (section 3.1) and the way a [2D image  $\rightarrow$  embeddings] estimator is trained (section 3.2). The method overview is illustrated in Figure 2.

#### 3.1 DENSEMARKS SPACE

The architecture of our pipeline consists of two key components: the canonical space where the embeddings reside, and the embedder, the task of which is to map an image into this space. The requirements that we set for the space are: (1) interpretable and queryable (the user can query a point in the space by looking at a typical arrangement of regions in it); (2) structured (regions are meaningful and don't overlap); (3) complete (contains the whole head, including the parts that are not trivial to annotate, such as hair and accessories); (4) smooth and continuous (images will be mapped to a continuous manifold of the space, with regions not getting abrupt and intersecting each other).

Additionally, we know that human heads are 3D objects. Even though UV (i.e., 2D canonical space) is a typical surface representation for heads, it's not the most precise representation due to modeling a complete head, including, e.g., hair and accessories, being not trivial in UV space and featuring seams (Ianina et al., 2022). Because of this, we decide to represent our canonical space as a unit cube in 3D and make the canonical embeddings the locations in this cube.

The interpretability requirement (1) and structure requirement (2) are enforced via landmark and segmentation losses, defined further in section 3.2.

The completeness requirement (3) is enforced with the way the embedder is supervised (also see section 3.2). For that purpose, we add a latent grid on top of the cube of a given resolution  $N_d \times N_d \times N_d$  and attach a  $D$ -dimensional latent feature to each element of the voxel grid, thus forming a learnable matrix  $E_{\text{raw}} \in \mathbb{R}^{(N_d)^3 \times D}$ . Each latent feature contains a highly-dimensional info about the given location in the canonical space.

Finally, to promote the smoothness requirement (4), we apply spatial smoothness to the matrix  $E_{\text{raw}}$  via a 3D Gaussian filter with a strength of  $\sigma$ , thus creating a latent feature grid  $E = \text{gaussian\_filter\_3D}(E_{\text{raw}}, \sigma)$ . This encourages the predicted embeddings from the embedder to be smoother, since the semantics of the close points in the cube will be similar and smoothly changing.

Note that the use of matrix  $E$  is inspired by the similar matrix of latent features used in functional maps, e.g., in CSE (Neverova et al., 2020), that is typically smoothed via a Laplace-Beltrami operator (Lévy, 2006). From a different standpoint, the operation of querying the space can also be seen as an attention operation, where the locations are queries (same as keys in this context) and the latent grid features are values. By aggregating the values at real-valued query locations with trilinear interpolation weights, we obtain the resulting semantic features at a given location.

### 3.2 EMBEDDER TRAINING

Our goal is to learn a monocular embedder  $\psi_\theta : I \rightarrow I_C$ , where  $I \in \mathbb{R}^{H \times W \times 3}$  is an input RGB image and  $I_C \in \mathbb{R}^{H \times W \times 3}$  is the predicted canonical embeddings for each pixel.

The network consists of a Vision Transformer backbone that predicts a feature map, which is further gradually upscaled through a sequence of convolutional layers to match the input resolution.

To train this network, at each training step, we pass two input images  $I^1, I^2 \in \mathbb{R}^{H \times W \times 3}$  through the embedder  $\psi_\theta$  and obtain corresponding predictions  $I_C^1 = \psi_\theta(I_1), I_C^2 = \psi_\theta(I_2)$ , both in  $\mathbb{R}^{H \times W \times 3}$ . For these two images, we assume having a number of ground truth pixel correspondences between them  $(K_{\text{gt}}^1, K_{\text{gt}}^2) = (\{(i_1^1, j_1^1), \dots, (i_P^1, j_P^1)\}, \{(i_1^2, j_1^2), \dots, (i_P^2, j_P^2)\})$ . These correspondences could be coming from any off-the-shelf pairwise matching algorithm. In our case, we obtain them from a point tracker inferred over individual talking head videos, as we found best in practice. Because of this, in our training procedure, images  $I_1$  and  $I_2$  are always coming from the same talking head video, but can represent arbitrarily close or far frames of the same video.

Embeddings  $I_C^1 = \psi_\theta(I_1)$  and  $I_C^2 = \psi_\theta(I_2)$  point to some real-valued locations in the canonical space. For each of those, we extract their corresponding  $D$ -dimensional semantic features via trilinear interpolation (Trilerp) (Bourke, 1999):  $I_{\text{feat}}^1 \in \mathbb{R}^{H \times W \times D}, I_{\text{feat}}^2 \in \mathbb{R}^{H \times W \times D}$ , where  $(I_{\text{feat}}^1)_{ij} = \text{Trilerp}(E, (I_C^1)_{ij}), (I_{\text{feat}}^2)_{ij} = \text{Trilerp}(E, (I_C^2)_{ij})$ .

In order to supervise our network, we encourage the features  $I_{\text{feat}}^1, I_{\text{feat}}^2$  to be close at the positions, defined by ground truth correspondences  $(K_{\text{gt}}^1, K_{\text{gt}}^2)$ , and far for other pairs of points. More formally, we first extract semantic features at the integer spatial positions of the ground truth correspondences, yielding tensors of queried features  $\text{Feat}^1, \text{Feat}^2 \in \mathbb{R}^{P \times D}$ ,  $\text{Feat}_p^1 = I_{\text{feat}}^1[(K_{\text{gt}}^1)_p], \text{Feat}_p^2 = I_{\text{feat}}^2[(K_{\text{gt}}^2)_p]$ . To promote the corresponding features of the first and second image to be close (*positive pairs*) and the others to be far (*negative pairs*), we construct a contrastive loss similar to CLIP Loss (Radford et al., 2021) that requires the pairwise matrix of cosine distances to be close to an identity matrix:

$$\mathcal{L}_{\theta, E}^{\text{contr}}(\text{Feat}^1, \text{Feat}^2) = \|(\text{norm}(\text{Feat}^1))(\text{norm}(\text{Feat}^1))^T - I\|_F,$$

where *norm* is a row-wise normalization operation.

Additionally, we apply a number of regularizations. To reduce ambiguity of the learned canonical space, we impose the locations of standard 300W Sagonas et al. (2013) format face landmarks to be close to the predefined locations in the cube. This is implemented via inferring an off-the-shelf landmark predictor on images  $I_1, I_2$ , thus obtaining ground truth landmark locations  $(l_1^1, \dots, l_{68}^1), (l_1^2, \dots, l_{68}^2)$ , and anchoring them to the predefined locations  $L_k \in \mathbb{R}^3, k = 1, \dots, 68$  in the unit cube:

$$\mathcal{L}_\theta^{\text{lmks}}(I_C | l) = \sum_{k=1}^{68} |I_C^1[l_k] - L_k|$$

To further correlate the predicted canonical embeddings with image semantics, we add a trainable segmentation head  $S_\xi$ , consisting of a single conv1x1 layer. For each of the images, this head receives the extracted semantic features (either  $\text{Feat}^1$  or  $\text{Feat}^2$ ) and returns the predicted logits of probabilities of class regions (face parsing) – either  $S^1 = S_\xi(\text{Feat}^1)$ , or  $S^2 = S_\xi(\text{Feat}^2)$ , both in  $\mathbb{R}^{H \times W \times N_S}$ . The segmentation loss expression compares each of the predicted masks  $S \in \{S^1, S^2\}$

270 to the corresponding ground truth mask  $S_{\text{gt}} \in \mathbb{R}^{H \times W \times N_S}$ , obtained by an off-the-shelf face parser:  
 271

$$272 \quad l^{\text{segm}}(S | S_{\text{gt}}) = \sum_{i,j} \text{cross\_entropy}(S[i, j], S_{\text{gt}}[i, j]) \\ 273 \\ 274$$

275 The overall loss is as follows:  
 276

$$277 \quad \mathcal{L}_{\theta, E, \xi}(\cdot) = \mathcal{L}_{\theta, E}^{\text{contr}}(\text{Feat}^1, \text{Feat}^2) \\ 278 \quad + \lambda_{\text{lmks}}(l_{\theta}^{\text{lmks}}(\mathbf{I}_C^1 | \mathbf{l}^1) + l_{\theta}^{\text{lmks}}(\mathbf{I}_C^2 | \mathbf{l}^2)) \\ 279 \quad + \lambda_{\text{segm}}(l^{\text{segm}}(S^1 | S_{\text{gt}}^1) + l^{\text{segm}}(S^2 | S_{\text{gt}}^2)) \\ 280$$

## 281 4 EXPERIMENTS

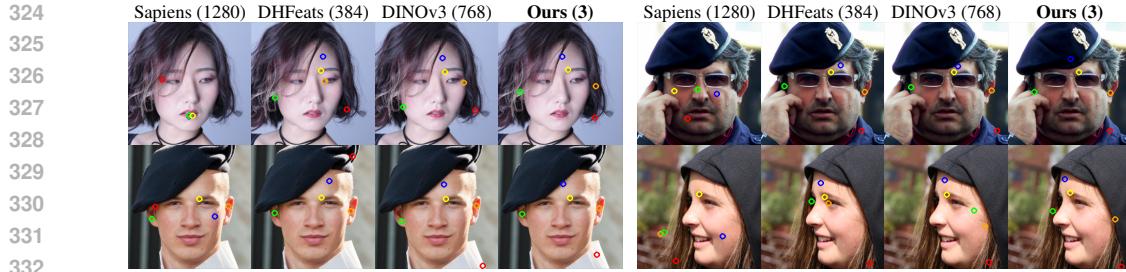
### 283 4.1 EXPERIMENTAL SETUP

285 **Data.** We train our method on CelebV-HQ dataset (Zhu et al., 2022) of 35K in-the-wild talking  
 286 head videos of interview style. To obtain ground truth correspondences  $(K_{\text{gt}}^1, K_{\text{gt}}^2)$ , we run Co-  
 287 Tracker3 (Karaev et al., 2024a) on these videos. As an input set of points to track, we take the whole  
 288 foreground region of the first frame (estimated by GroundedSAM2 (Ren et al., 2024) prompted with  
 289 the text “*person*”) and sample points uniformly in that region (see an example in Fig. 1 (*left*)).  
 290 Videos were discarded if there were either too few tracks found (fewer than 80) or foreground seg-  
 291 mentation failed, resulting in 32K videos left. The number of point tracks found did not exceed  
 292 400. 100 randomly sampled videos have been held out for the evaluation and used in the results  
 293 described below. Each training batch is formed by uniformly sampling two random frames from a  
 294 sample video from the constructed annotated dataset. All videos are resized to the (512, 512) res-  
 295 olution in advance and fed to the embedder in that resolution. For augmentation, we use random  
 296 shift (in [-10%, 10%] range), scale ([-10%, 10%]), and rotation ([-18°, 18°]), each with a chance of  
 297 50%. Points which are no longer visible after augmentation are no longer accounted in training. For  
 298 the landmark loss, we extract 70 manually selected landmarks (full face border, landmarks on eyes,  
 299 nose, and mouth) via Mediapipe (Lugaresi et al., 2019b). Ground truth segmentation masks are ob-  
 300 tained via FaRL (Zheng et al., 2022) and are further refined on the borders via face-parsing (Jonathan  
 301 Dinu, 2025; Xie et al., 2021), which works better in practice on non-face regions of the head.

302 **Architecture and training.** To make use of strong pretraining, we initialize the embedder with a  
 303 pre-trained DINOv3 (Siméoni et al., 2025) checkpoint and add DPT head (Ranftl et al., 2021) to  
 304 output an image of the same spatial resolution as the input ( $512 \times 512$ ). Matrix  $E$  is initialized from  
 305 a Gaussian distribution  $\mathcal{N}(0, 1)$ . We use  $\lambda_{\text{segm}} = 1$  for the segmentation loss and  $\lambda_{\text{lmks}} = 50$  for  
 306 the landmark loss. For optimization, we employ the AdamW (Loshchilov, 2017) optimizer with a  
 307 learning rate  $5 \cdot 10^{-5}$  for the backbone of  $\phi_{\theta}$ , learning rate of  $10^{-4}$  for DPT head, and  $10^{-3}$  for the  
 308 latent features  $E$ . The schedule for all learning rates was cosine annealing with an overall number  
 309 of steps of 140K and a warmup for 2'800 steps. Weight decay of  $10^{-4}$  was applied to the network  
 310 parameters  $\theta$  and  $\xi$ , except for normalization layers. The whole pipeline is trained for 140k training  
 311 steps using 8 pairs of images per batch on a single NVIDIA RTX 3090 Ti GPU for 1.5 days.

### 312 4.2 RESULTS

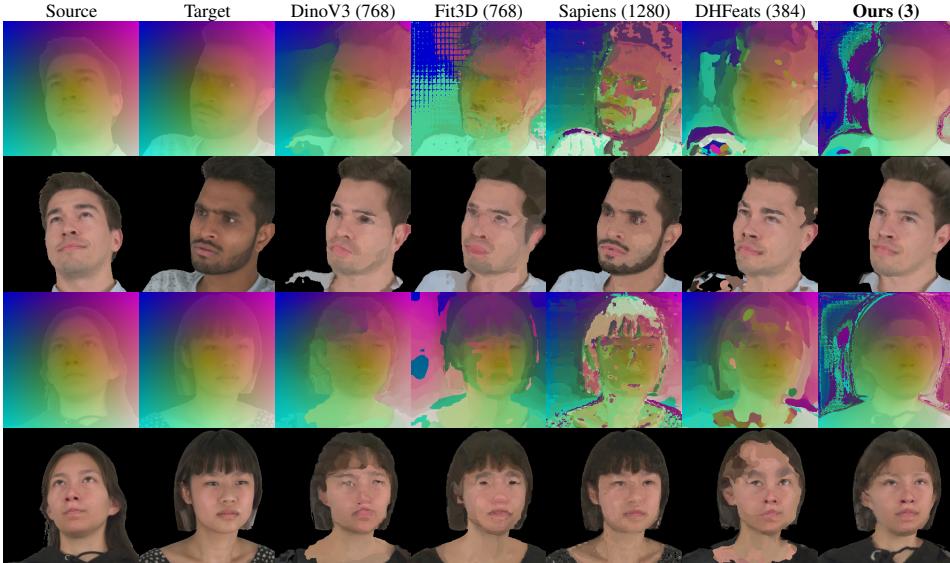
313 **Point querying.** The requirement of the canonical space is that the same semantic points will  
 314 have a fixed location in the cube, regardless of the person’s identity. We test this on a number of  
 315 points that have distinct semantics: points on hair, ear centers, forehead center, eyebrow corners.  
 316 To find each semantic point, we manually annotated 7 sample images from CelebV-HQ, inferred  
 317 the trained embedder, and averaged predicted locations in the cube for each annotated point. We  
 318 use the obtained location as a reference to find the nearest neighbor in the other image among their  
 319 predicted embeddings. Results are demonstrated in Fig. 3. There, we compare against state-of-  
 320 the-art dense feature extractors, the embeddings of which provide rich semantic information for a  
 321 neighbor search: DINOv3 (Siméoni et al., 2025) (embedding dimension: 768), Sapiens (Khirodkar  
 322 et al., 2024) (1280), Diffusion Hyperfeatures (Luo et al., 2023) (384), Fit3D (Yue et al., 2024) (768).  
 323 For these methods, semantic points are also estimated by averaging predicted embeddings. Despite  
 324 using a significantly smaller vector dimension (3) to store semantics in the embedding, our method



333      Figure 3: Point querying. We select a specific point on a few images and find the reference embedding  
 334      by averaging the embeddings predicted by each of the models in its location. Points: **red** = on  
 335      the left side of long hair region, **green** = center of the right ear, **orange** = center of the left ear, **blue**  
 336      = forehead center, **yellow** = left eyebrow corner. We indicate the embedding dimension in brackets.  
 337



346      Figure 4: Semantic regions on head images can be located via selecting corresponding volumetric  
 347      regions in the canonical space. Blue: forehead center, green and orange: ears, yellow: skin near the  
 348      left eyebrow corner.  
 349



369      Figure 5: Dense warping. Here, we copy pixels from source to target based on the target  $\rightarrow$  source  
 370      nearest neighbors search in the space of embeddings, predicted by each model (*even rows*). For  
 371      clarity, mapping of meshgrid-like coordinates, blended with RGB, is shown additionally (*odd rows*).  
 372      Even though deep feature extractors provide valuable matches, they are either matching colors, not  
 373      semantics (Sapiens (Khirodkar et al., 2024), DHFeats (Luo et al., 2023)), or feature significant artifacts  
 374      (DinoV3 (Siméoni et al., 2025), Fit3D (Yue et al., 2024)), thus being less reliable for matching.  
 375  
 376  
 377

can find a corresponding region for challenging views better. Note that our method is also robust to strong face or head occlusions.

378  
 379 Table 1: Quantitative comparison. On same-person pairs of images from Nerssemble (Kirschstein  
 380 et al., 2023), we evaluate the quality of correspondences that arise from matching nearest neighbor  
 381 embeddings. Similarly, on cross-person pairs, we evaluate the consistency and identity preservation.  
 382  
 383

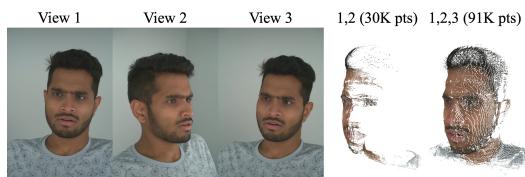
	Same-person		Cross-person	
	Matching quality MAE ↓	RMSE ↓	ArcFace ↑	Met3R ↓
DINOv3 (Siméoni et al., 2025)	7.6	12.69	0.266	0.460
Fit3D (Yue et al., 2024)	12.75	21.83	0.236	0.558
Hyperfeatures (Luo et al., 2023)	8.26	13.29	0.329	0.454
Sapiens (Khirodkar et al., 2024)	14.88	24.12	0.167	0.595
<b>Ours</b>	<b>3.68</b>	<b>5.9</b>	<b>0.384</b>	<b>0.388</b>



410 Figure 6: Monocular tracking. We evaluate our method on downstream application of applying  
 411 a state-of-the-art off-the-shelf head tracker (Qian, 2024) to track a 3D Morphable Model template  
 412 (FLAME (Li et al., 2017a)) over a monocular video. By default, this tracker relies on standard 68  
 413 face landmarks and photometric loss. Estimating a DenseMarks texture of FLAME and applying an  
 414 additional photometric loss to match it with estimated embeddings greatly improves the robustness  
 415 of the tracker, especially for extreme poses.

416  
 417 **Region selection.** In Fig. 4, we demonstrate  
 418 how the same volumetric region in the canoni-  
 419 cal space is mapped onto images of people. The  
 420 regions are initially selected on 7 random im-  
 421 ages manually and averaged (via a voting pro-  
 422 cedure) in the cube space.

423  
 424 **Dense warping.** To demonstrate the seman-  
 425 tic consistency of embeddings predicted for the  
 426 whole image, not only specific points or re-  
 427 gions, we demonstrate the warping by embed-  
 428 dings in Fig. 5, evaluated on pairs of different people from the Nerssemble dataset (Kirschstein et al.,  
 429 2023). For each target image pixel, we replace its color with the color of the nearest neighbor by  
 430 embedding in the source. We expect the warping to be semantically meaningful and smooth. It is  
 431 observed that when we match nearest neighbors by Diffusion Hyperfeatures and especially Sapiens  
 432 embeddings, the matches turn out to be based on the color similarity, not the semantic similarity.  
 433 DINOv3 and Fit3D appear more semantically meaningful but often feature artifacts, making the cor-  
 434 respondences imprecise, as best observed in the mapping rows in the figure. To evaluate the quality



435 Figure 7: Stereo Reconstruction. We triangulate  
 436 2-view and 3-view correspondences of our  
 437 representations using known camera parameters in  
 438 Nerssemble (Kirschstein et al., 2023).

432 of the mapping, we estimate face recognition similarity based on ArcFace (Deng et al., 2019) be-  
 433 tween the source image and the mapping result, as well as the view-consistency metric Met3R (Asim  
 434 et al., 2025), and show the results in Table 1.

435 **Geometric consistency.** To assess qualitatively and quantitatively the precision of the estimated cor-  
 436 respondences through our embeddings, we repeat the Dense Warping experiment in a similar way  
 437 for the (source, target) pairs of images of the same person, not different people, repeated over various  
 438 people from the Nerensemble dataset. In Table 1, we demonstrate the evaluation of the correctness of  
 439 the estimated correspondences between source and target, averaged over ten people from Nerensem-  
 440 ble. As a source of ground truth correspondences, we estimate a complete head mesh from all 16  
 441 cameras via GS2Mesh (Wolf et al., 2024) and sample 1K random mesh vertices. The embeddings  
 442 are evaluated in the projected locations of these vertices.

443 **Losses ablation.** Even though the network can  
 444 learn without introduced constraints on land-  
 445 mark locations in the cube and segmentation  
 446 loss, we demonstrate that the finding char-  
 447 acteristic points and regions becomes more  
 448 problematic in Fig. 8. This is explained by a less  
 449 semantically constrained canonical space.

450 **Monocular tracking.** As an example  
 451 application of our method, we take a  
 452 highly-performing off-the-shelf head tracker,  
 453 VHAP (Qian, 2024), which supports estima-  
 454 tion of the FLAME parametric head model (Li  
 455 et al., 2017a). It relies on a standard 300-W  
 456 set of 68 sparse landmarks (Sagonas et al.,  
 457 2013) for rigid alignment of the template and  
 458 optimizes for the shape, pose, and expression  
 459 parameters of FLAME, through estimating  
 460 RGB texture in the FLAME UV space and  
 461 applying photometric loss. Even though VHAP  
 462 excels in multi-view settings, monocular  
 463 videos can remain challenging due to poten-  
 464 tially failing landmark detection, occlusions,  
 465 and extreme viewpoints. To aid the tracker in  
 466 these situations, we add another photometric  
 467 loss that is based on estimating a 3-dimensional  
 468 UV texture of DenseMarks embeddings that is  
 469 compared to the embeddings predicted by the trained  
 470 embedder for each video frame independently.  
 471 We run tracking on in-the-wild monocular videos  
 472 with different challenging conditions such as strong/fast head rotation, severe hair/accessories  
 473 occlusions, very close/far cameras. The results are demonstrated in Figure 6. Our method improves  
 474 robustness the most in cases of extreme poses and yields better alignment in challenging regions,  
 475 such as neck and ears. We demonstrate the results of tracking over the complete videos in the  
 476 Supplementary Video.

477 **Stereo Reconstruction.** In Fig. 7, we demon-  
 478 strate that triangulating 2+ images can be done purely  
 479 using embeddings from our model, on the example of a sample from Nerensemble with known camera  
 480 poses and intrinsics. This way, we demon-  
 481 strate the capabilities of [multi-view]-stereo and dense  
 482 estimation.

## 483 5 CONCLUSION

484 We propose a novel representation for human head images and an embedder for dense estimation.  
 485 The resulting low-dimensional (3D) embeddings are consistent across views and subjects, enabling  
 486 reliable matching of challenging regions like hair. Despite their compactness, they outperform high-  
 487 dimensional features from foundation models in geometry-aware tasks like tracking, while benefit-  
 488 ing from VFM pretraining. Future work could extend our approach to full bodies and other domains,  
 489 which would be anticipated with the appearance of publicly available high-resolution data collec-  
 490 tions.



491 Figure 8: Removing the landmark or segmen-  
 492 tation loss makes region finding much less reliable.  
 493 **Blue:** forehead center, **green** and **orange**: ears,  
 494 **yellow**: skin near the left eyebrow corner.

486 6 APPENDIX  
487488 **Use of LLMs.** We used LLMs for expanding our knowledge regarding the latest related work.  
489490 REFERENCES  
491

492 Hervé Abdi and Lynne J Williams. Principal component analysis. *Wiley interdisciplinary reviews: 493 computational statistics*, 2(4):433–459, 2010.  
494

495 Mohammad Asim, Christopher Wewer, Thomas Wimmer, Bernt Schiele, and Jan Eric Lenssen.  
496 Met3r: Measuring multi-view consistency in generated images. In *Proceedings of the Computer  
497 Vision and Pattern Recognition Conference*, pp. 6034–6044, 2025.

498 Volker Blanz and Thomas Vetter. A morphable model for the synthesis of 3d faces. In *Proceedings  
499 of the 26th Annual Conference on Computer Graphics and Interactive Techniques*, SIGGRAPH  
500 '99, pp. 187–194, USA, 1999. ACM Press/Addison-Wesley Publishing Co. ISBN 0201485605.  
501 doi: 10.1145/311535.311556. URL <https://doi.org/10.1145/311535.311556>.

502 Federica Bogo, Angjoo Kanazawa, Christoph Lassner, Peter Gehler, Javier Romero, and Michael J  
503 Black. Keep it smpl: Automatic estimation of 3d human pose and shape from a single image. In  
504 *European conference on computer vision*, pp. 561–578. Springer, 2016.  
505

506 Paul Bourke. Interpolation methods. *Miscellaneous: projection, modelling, rendering*, 1(10), 1999.  
507

508 Adrian Bulat and Georgios Tzimiropoulos. How far are we from solving the 2d & 3d face align-  
509 ment problem? (and a dataset of 230,000 3d facial landmarks). In *International Conference on  
510 Computer Vision*, 2017.

511 Dongliang Cao and Florian Bernard. Unsupervised deep multi-shape matching. In *European con-  
512 ference on computer vision*, pp. 55–71. Springer, 2022.  
513

514 Dongliang Cao, Paul Roetzer, and Florian Bernard. Unsupervised learning of robust spectral shape  
515 matching. *arXiv preprint arXiv:2304.14419*, 2023.

516 Z. Cao, G. Hidalgo Martinez, T. Simon, S. Wei, and Y. A. Sheikh. Openpose: Realtime multi-person  
517 2d pose estimation using part affinity fields. *IEEE Transactions on Pattern Analysis and Machine  
518 Intelligence*, 2019.

519 Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and  
520 Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of  
521 the International Conference on Computer Vision (ICCV)*, 2021.

523 Eric R Chan, Connor Z Lin, Matthew A Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio  
524 Gallo, Leonidas J Guibas, Jonathan Tremblay, Sameh Khamis, et al. Efficient geometry-aware 3d  
525 generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision  
526 and pattern recognition*, pp. 16123–16133, 2022.

527 Xingyu Chen, Yue Chen, Yuliang Xiu, Andreas Geiger, and Anpei Chen. Easi3r: Estimating disen-  
528 tangled motion from dust3r without training. *arXiv preprint arXiv:2503.24391*, 2025a.  
529

530 Zhuoguang Chen, Minghui Qin, Tianyuan Yuan, Zhe Liu, and Hang Zhao. Long3r: Long sequence  
531 streaming 3d reconstruction. *arXiv preprint arXiv:2507.18255*, 2025b.  
532

533 Seokju Cho, Jiahui Huang, Jisu Nam, Honggyu An, Seungryong Kim, and Joon-Young Lee. Local  
534 all-pair correspondence for point tracking. *arXiv preprint arXiv:2407.15420*, 2024.

535 Hang Dai, Nick Pears, William Smith, and Christian Duncan. Statistical modeling of craniofacial  
536 shape and texture. *International Journal of Computer Vision*, 128(2):547–571, 2020.  
537

538 Radek Daněček, Michael J Black, and Timo Bolkart. Emoca: Emotion driven monocular face  
539 capture and animation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and  
Pattern Recognition*, pp. 20311–20322, 2022.

540 Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. Arcface: Additive angular margin  
 541 loss for deep face recognition. In *Proceedings of the IEEE/CVF conference on computer vision*  
 542 and pattern recognition, pp. 4690–4699, 2019.

543 Carl Doersch, Ankush Gupta, Larisa Markeeva, Adria Recasens, Lucas Smaira, Yusuf Aytar, Joao  
 544 Carreira, Andrew Zisserman, and Yi Yang. Tap-vid: A benchmark for tracking any point in a  
 545 video. *Advances in Neural Information Processing Systems*, 35:13610–13626, 2022.

546 Carl Doersch, Yi Yang, Mel Vecerik, Dilara Gokay, Ankush Gupta, Yusuf Aytar, Joao Carreira,  
 547 and Andrew Zisserman. Tapir: Tracking any point with per-frame initialization and temporal  
 548 refinement. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp.  
 549 10061–10072, 2023.

550 Carl Doersch, Pauline Luc, Yi Yang, Dilara Gokay, Skanda Koppula, Ankush Gupta, Joseph Hey-  
 551 ward, Ignacio Rocco, Ross Goroshin, Joao Carreira, et al. Bootstrap: Bootstrapped training for  
 552 tracking-any-point. In *Proceedings of the Asian Conference on Computer Vision*, pp. 3257–3274,  
 553 2024.

554 Zijian Dong, Xu Chen, Jinlong Yang, Michael J Black, Otmar Hilliges, and Andreas Geiger. Ag3d:  
 555 Learning to generate 3d avatars from 2d image collections. In *Proceedings of the IEEE/CVF*  
 556 *international conference on computer vision*, pp. 14916–14927, 2023.

557 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas  
 558 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An  
 559 image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint*  
 560 *arXiv:2010.11929*, 2020.

561 Niladri Shekhar Dutt, Sanjeev Muralikrishnan, and Niloy J Mitra. Diffusion 3d features (diff3f):  
 562 Decorating untextured shapes with distilled semantic features. In *Proceedings of the IEEE/CVF*  
 563 *Conference on Computer Vision and Pattern Recognition*, pp. 4494–4504, 2024.

564 Ariel Ephrat, Inbar Mosseri, Oran Lang, Tali Dekel, Kevin Wilson, Avinatan Hassidim, William T  
 565 Freeman, and Michael Rubinstein. Looking to listen at the cocktail party: A speaker-independent  
 566 audio-visual model for speech separation. *arXiv preprint arXiv:1804.03619*, 2018.

567 Haiwen Feng, Junyi Zhang, Qianqian Wang, Yufei Ye, Pengcheng Yu, Michael J Black, Trevor  
 568 Darrell, and Angjoo Kanazawa. St4rtrack: Simultaneous 4d reconstruction and tracking in the  
 569 world. *arXiv preprint arXiv:2504.13152*, 2025.

570 Yao Feng, Fan Wu, Xiaohu Shao, Yanfeng Wang, and Xi Zhou. Joint 3d face reconstruction and  
 571 dense alignment with position map regression network. In *Proceedings of the European confer-  
 572 ence on computer vision (ECCV)*, pp. 534–551, 2018.

573 Yao Feng, Haiwen Feng, Michael J Black, and Timo Bolkart. Learning an animatable detailed 3d  
 574 face model from in-the-wild images. *ACM Transactions on Graphics (ToG)*, 40(4):1–13, 2021.

575 Guy Gafni, Justus Thies, Michael Zollhofer, and Matthias Nießner. Dynamic neural radiance fields  
 576 for monocular 4d facial avatar reconstruction. In *Proceedings of the IEEE/CVF Conference on*  
 577 *Computer Vision and Pattern Recognition*, pp. 8649–8658, 2021.

578 Simon Giebenhain, Tobias Kirschstein, Markos Georgopoulos, Martin Rünz, Lourdes Agapito, and  
 579 Matthias Nießner. Mononphm: Dynamic head reconstruction from monocular videos. In *Proc.  
 580 IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2024.

581 Simon Giebenhain, Tobias Kirschstein, Martin Rünz, Lourdes Agapito, and Matthias Nießner.  
 582 Pixel3dmm: Versatile screen-space priors for single-image 3d face reconstruction. *arXiv preprint*  
 583 *arXiv:2505.00615*, 2025.

584 Philip-William Grassal, Malte Prinzler, Titus Leistner, Carsten Rother, Matthias Nießner, and Justus  
 585 Thies. Neural head avatars from monocular rgb videos. *arXiv preprint arXiv:2112.01554*, 2021.

586 Rıza Alp Güler, Natalia Neverova, and Iasonas Kokkinos. Densepose: Dense human pose estimation  
 587 in the wild. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,  
 588 pp. 7297–7306, 2018.

594 Oshri Halimi, Or Litany, Emanuele Rodola, Alex M Bronstein, and Ron Kimmel. Unsupervised  
 595 learning of dense shape correspondence. In *Proceedings of the IEEE/CVF Conference on Com-*  
 596 *puter Vision and Pattern Recognition*, pp. 4370–4379, 2019.

597

598 Kai Han, Yunhe Wang, Hanting Chen, Xinghao Chen, Jianyuan Guo, Zhenhua Liu, Yehui Tang,  
 599 An Xiao, Chunjing Xu, Yixing Xu, et al. A survey on vision transformer. *IEEE transactions on*  
 600 *pattern analysis and machine intelligence*, 45(1):87–110, 2022.

601 Adam W Harley, Zhaoyuan Fang, and Katerina Fragkiadaki. Particle video revisited: Tracking  
 602 through occlusions using point trajectories. In *European Conference on Computer Vision*, pp.  
 603 59–75. Springer, 2022.

604 James V Haxby, Elizabeth A Hoffman, and M Ida Gobbini. The distributed human neural system  
 605 for face perception. *Trends in cognitive sciences*, 4(6):223–233, 2000.

606

607 Eric Hedlin, Gopal Sharma, Shweta Mahajan, Hossam Isack, Abhishek Kar, Andrea Tagliasacchi,  
 608 and Kwang Moo Yi. Unsupervised semantic correspondence using stable diffusion. *Advances in*  
 609 *Neural Information Processing Systems*, 36:8266–8279, 2023.

610 Yinghao Huang, Omid Taheri, Michael J Black, and Dimitrios Tzionas. Intercap: Joint markerless  
 611 3d tracking of humans and objects in interaction. In *DAGM German Conference on Pattern*  
 612 *Recognition*, pp. 281–299. Springer, 2022.

613

614 Anastasia Ianina, Nikolaos Sarafianos, Yuanlu Xu, Ignacio Rocco, and Tony Tung. Bodymap:  
 615 Learning full-body dense correspondence map. In *Proceedings of the IEEE/CVF conference on*  
 616 *computer vision and pattern recognition*, pp. 13286–13295, 2022.

617 Jiaxi Jiang, Paul Streli, Huajian Qiu, Andreas Fender, Larissa Laich, Patrick Snape, and Christian  
 618 Holz. Avataposer: Articulated full-body pose tracking from sparse motion sensing. In *European*  
 619 *conference on computer vision*, pp. 443–460. Springer, 2022.

620

621 Sheng Jin, Lumin Xu, Jin Xu, Can Wang, Wentao Liu, Chen Qian, Wanli Ouyang, and Ping Luo.  
 622 Whole-body human pose estimation in the wild. In *European Conference on Computer Vision*,  
 623 pp. 196–214. Springer, 2020.

624

625 Ian Jolliffe. Principal component analysis. In *International encyclopedia of statistical science*, pp.  
 1094–1096. Springer, 2011.

626

627 Jonathan Dinu. *jonathandinu/face-parsing*: Face parsing model (fine-tuned from segformer on  
 628 celebamask-hq). <https://huggingface.co/jonathandinu/face-parsing>, 2025.  
 629 Accessed: 2025-09-25.

630

631 Angjoo Kanazawa, Michael J Black, David W Jacobs, and Jitendra Malik. End-to-end recovery of  
 632 human shape and pose. In *Proceedings of the IEEE conference on computer vision and pattern*  
 633 *recognition*, pp. 7122–7131, 2018.

634

635 Nikita Karaev, Iurii Makarov, Jianyuan Wang, Natalia Neverova, Andrea Vedaldi, and Christian  
 636 Rupprecht. Cotracker3: Simpler and better point tracking by pseudo-labelling real videos. *arXiv*  
 637 *preprint arXiv:2410.11831*, 2024a.

638

639 Nikita Karaev, Ignacio Rocco, Benjamin Graham, Natalia Neverova, Andrea Vedaldi, and Christian  
 640 Rupprecht. Cotracker: It is better to track together. In *European conference on computer vision*,  
 641 pp. 18–35. Springer, 2024b.

642

643 Yoni Kasten, Wuyue Lu, and Haggai Maron. Fast encoder-based 3d from casual videos via point  
 644 track processing. *Advances in Neural Information Processing Systems*, 37:96150–96180, 2024.

645

646 Rawal Khirodkar, Timur Bagautdinov, Julieta Martinez, Su Zhaoen, Austin James, Peter Selednik,  
 647 Stuart Anderson, and Shunsuke Saito. Sapiens: Foundation for human vision models. In *Euro-*  
 648 *pean Conference on Computer Vision*, pp. 206–228. Springer, 2024.

649

650 Inès Hyeonsu Kim, Seokju Cho, Jahyeok Koo, Junghyun Park, Jiahui Huang, Joon-Young Lee,  
 651 and Seungryong Kim. Learning to track any points from human motion. *arXiv preprint*  
 652 *arXiv:2507.06233*, 2025.

648 Tobias Kirschstein, Shenhan Qian, Simon Giebenhain, Tim Walter, and Matthias Nießner. Nerseme-  
 649 ble: Multi-view radiance field reconstruction of human heads. *ACM Transactions on Graphics*  
 650 (*TOG*), 42(4):1–14, 2023.

651 Muhammed Kocabas, Nikos Athanasiou, and Michael J Black. Vibe: Video inference for human  
 652 body pose and shape estimation. In *Proceedings of the IEEE/CVF conference on computer vision*  
 653 and *pattern recognition*, pp. 5253–5263, 2020.

654 Akshay Krishnan, Abhijit Kundu, Kevins-Kokitsi Maninis, James Hays, and Matthew Brown.  
 655 Omnimoces: A unified noces dataset and model for 3d lifting of 2d objects. In *European Conference*  
 656 on *Computer Vision*, pp. 127–145. Springer, 2024.

657 Vincent Leroy, Yohann Cabon, and Jérôme Revaud. Grounding image matching in 3d with mast3r.  
 658 In *European Conference on Computer Vision*, pp. 71–91. Springer, 2024.

659 Bruno Lévy. Laplace-beltrami eigenfunctions towards an algorithm that “understands” geometry. In  
 660 *IEEE International Conference on Shape Modeling and Applications 2006 (SMI’06)*, pp. 13–13.  
 661 IEEE, 2006.

662 Hongyang Li, Hao Zhang, Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, and Lei Zhang. Taptr:  
 663 Tracking any point with transformers as detection. In *European Conference on Computer Vision*,  
 664 pp. 57–75. Springer, 2024.

665 Hui Li, Zidong Guo, Seon-Min Rhee, Seungju Han, and Jae-Joon Han. Towards accurate facial  
 666 landmark detection via cascaded transformers. In *Proceedings of the IEEE/CVF conference on*  
 667 *computer vision and pattern recognition*, pp. 4176–4185, 2022.

668 Tianye Li, Timo Bolkart, Michael. J. Black, Hao Li, and Javier Romero. Learning a model of facial  
 669 shape and expression from 4D scans. *ACM Transactions on Graphics, (Proc. SIGGRAPH Asia)*,  
 670 36(6):194:1–194:17, 2017a. URL <https://doi.org/10.1145/3130800.3130813>.

671 Tianye Li, Timo Bolkart, Michael. J. Black, Hao Li, and Javier Romero. Learning a model of facial  
 672 shape and expression from 4D scans. *ACM Transactions on Graphics, (Proc. SIGGRAPH Asia)*,  
 673 36(6):194:1–194:17, 2017b. URL <https://doi.org/10.1145/3130800.3130813>.

674 Minghua Liu, Mikaela Angelina Uy, Donglai Xiang, Hao Su, Sanja Fidler, Nicholas Sharp, and  
 675 Jun Gao. Partfield: Learning 3d feature fields for part segmentation and beyond. *arXiv preprint*  
 676 *arXiv:2504.11451*, 2025.

677 Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black.  
 678 SMPL: A skinned multi-person linear model. *ACM Trans. Graphics (Proc. SIGGRAPH Asia)*,  
 679 34(6):248:1–248:16, October 2015.

680 I Loshchilov. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.

681 Camillo Lugaressi, Jiuqiang Tang, Hadon Nash, Chris McClanahan, Esha Ubweja, Michael Hays,  
 682 Fan Zhang, Chuo-Ling Chang, Ming Guang Yong, Juhyun Lee, et al. Mediapipe: A framework  
 683 for building perception pipelines. *arXiv preprint arXiv:1906.08172*, 2019a.

684 Camillo Lugaressi, Jiuqiang Tang, Hadon Nash, Chris McClanahan, Esha Ubweja, Michael Hays,  
 685 Fan Zhang, Chuo-Ling Chang, Ming Guang Yong, Juhyun Lee, et al. Mediapipe: A framework  
 686 for building perception pipelines. *arXiv preprint arXiv:1906.08172*, 2019b.

687 Jonathon Luiten, Georgios Kopanas, Bastian Leibe, and Deva Ramanan. Dynamic 3d gaussians:  
 688 Tracking by persistent dynamic view synthesis. In *2024 International Conference on 3D Vision*  
 689 (*3DV*), pp. 800–809. IEEE, 2024.

690 Grace Luo, Lisa Dunlap, Dong Huk Park, Aleksander Holynski, and Trevor Darrell. Diffusion  
 691 hyperfeatures: Searching through time and space for semantic correspondence. In *Advances in*  
 692 *Neural Information Processing Systems*, 2023.

693 Zhixiang Min, Bingbing Zhuang, Samuel Schulter, Buyu Liu, Enrique Dunn, and Manmohan Chan-  
 694 draker. Neurocs: Neural noces supervision for monocular 3d object localization. In *Proceedings of*  
 695 *the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 21404–21414, 2023.

702 Gyeongsik Moon, Shouu-I Yu, He Wen, Takaaki Shiratori, and Kyoung Mu Lee. Interhand2.6m: A  
 703 dataset and baseline for 3d interacting hand pose estimation from a single rgb image. In *European*  
 704 *Conference on Computer Vision*, pp. 548–564. Springer, 2020.

705 Natalia Neverova, David Novotny, Marc Szafraniec, Vasil Khalidov, Patrick Labatut, and Andrea  
 706 Vedaldi. Continuous surface embeddings. *Advances in Neural Information Processing Systems*,  
 707 33:17258–17270, 2020.

708 Maxime Oquab, Timothée Darcet, Theo Moutakanni, Huy V. Vo, Marc Szafraniec, Vasil Khalidov,  
 709 Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaeldin El-Nouby, Russell Howes, Po-Yao  
 710 Huang, Hu Xu, Vasu Sharma, Shang-Wen Li, Wojciech Galuba, Mike Rabbat, Mido Assran,  
 711 Nicolas Ballas, Gabriel Synnaeve, Ishan Misra, Herve Jegou, Julien Mairal, Patrick Labatut, Ar-  
 712 mand Joulin, and Piotr Bojanowski. Dinov2: Learning robust visual features without supervision,  
 713 2023.

714 Maks Ovsjanikov, Mirela Ben-Chen, Justin Solomon, Adrian Butscher, and Leonidas Guibas. Func-  
 715 tional maps: a flexible representation of maps between shapes. *ACM Transactions on Graphics*  
 716 (*ToG*), 31(4):1–11, 2012.

717 Keunhong Park, Utkarsh Sinha, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Steven M  
 718 Seitz, and Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. In *Proceedings of*  
 719 *the IEEE/CVF international conference on computer vision*, pp. 5865–5874, 2021.

720 Shenhan Qian. Vhap: Versatile head alignment with adaptive appearance priors, sep 2024.

721 Shenhan Qian, Tobias Kirschstein, Liam Schoneveld, Davide Davoli, Simon Giebenhain, and  
 722 Matthias Nießner. Gaussianavatars: Photorealistic head avatars with rigged 3d gaussians. In *Pro-  
 723 ceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 20299–  
 724 20309, 2024.

725 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,  
 726 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual  
 727 models from natural language supervision. In *International conference on machine learning*, pp.  
 728 8748–8763. PMLR, 2021.

729 René Ranftl, Alexey Bochkovskiy, and Vladlen Koltun. Vision transformers for dense prediction.  
 730 In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 12179–12188,  
 731 2021.

732 Tianhe Ren, Shilong Liu, Ailing Zeng, Jing Lin, Kunchang Li, He Cao, Jiayu Chen, Xinyu Huang,  
 733 Yukang Chen, Feng Yan, et al. Grounded sam: Assembling open-world models for diverse visual  
 734 tasks. *arXiv preprint arXiv:2401.14159*, 2024.

735 Emanuele Rodolà, Luca Cosmo, Michael M Bronstein, Andrea Torsello, and Daniel Cremers. Partial  
 736 functional correspondence. In *Computer graphics forum*, volume 36, pp. 222–236. Wiley Online  
 737 Library, 2017.

738 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
 739 resolution image synthesis with latent diffusion models, 2021.

740 Javier Romero, Dimitrios Tzionas, and Michael J. Black. Embodied hands: Modeling and capturing  
 741 hands and bodies together. *ACM Transactions on Graphics, (Proc. SIGGRAPH Asia)*, 36(6),  
 742 November 2017.

743 Christos Sagonas, Georgios Tzimiropoulos, Stefanos Zafeiriou, and Maja Pantic. 300 faces in-the-  
 744 wild challenge: The first facial landmark localization challenge. In *Proceedings of the IEEE*  
 745 *international conference on computer vision workshops*, pp. 397–403, 2013.

746 Soyong Shin, Juyong Kim, Eni Halilaj, and Michael J Black. Wham: Reconstructing world-  
 747 grounded humans with accurate 3d motion. In *Proceedings of the IEEE/CVF Conference on*  
 748 *Computer Vision and Pattern Recognition*, pp. 2070–2080, 2024.

756 Oriane Siméoni, Huy V. Vo, Maximilian Seitzer, Federico Baldassarre, Maxime Oquab, Cijo Jose,  
 757 Vasil Khalidov, Marc Szafraniec, Seungeun Yi, Michaël Ramamonjisoa, Francisco Massa, Daniel  
 758 Haziza, Luca Wehrstedt, Jianyuan Wang, Timothée Darct, Théo Moutakanni, Leonel Sentana,  
 759 Claire Roberts, Andrea Vedaldi, Jamie Tolan, John Brandt, Camille Couprie, Julien Mairal, Hervé  
 760 Jégou, Patrick Labatut, and Piotr Bojanowski. DINOv3, 2025. URL <https://arxiv.org/abs/2508.10104>.

762 Tomas Simon, Hanbyul Joo, Iain Matthews, and Yaser Sheikh. Hand keypoint detection in single  
 763 images using multiview bootstrapping. In *CVPR*, 2017.

764 Brandon Smart, Chuanxia Zheng, Iro Laina, and Victor Adrian Prisacariu. Splatt3r: Zero-shot  
 765 gaussian splatting from uncalibrated image pairs. *arXiv preprint arXiv:2408.13912*, 2024.

767 Edgar Sucar, Zihang Lai, Eldar Insafutdinov, and Andrea Vedaldi. Dynamic point maps: A versatile  
 768 representation for dynamic 3d reconstruction. *arXiv preprint arXiv:2503.16318*, 2025.

770 Felix Taubner, Prashant Raina, Mathieu Tuli, Eu Wern Teh, Chul Lee, and Jinmiao Huang. 3d face  
 771 tracking from 2d video through iterative dense uv to image flow. In *Proceedings of the IEEE/CVF*  
 772 *Conference on Computer Vision and Pattern Recognition*, pp. 1227–1237, 2024.

773 Justus Thies, Michael Zollhofer, Marc Stamminger, Christian Theobalt, and Matthias Nießner.  
 774 Face2face: Real-time face capture and reenactment of rgb videos. In *Proceedings of the IEEE*  
 775 *conference on computer vision and pattern recognition*, pp. 2387–2395, 2016.

776 Narek Tumanyan, Assaf Singer, Shai Bagon, and Tali Dekel. Dino-tracker: Taming dino for self-  
 777 supervised point tracking in a single video. In *European Conference on Computer Vision*, pp.  
 778 367–385. Springer, 2024.

779 He Wang, Srinath Sridhar, Jingwei Huang, Julien Valentin, Shuran Song, and Leonidas J Guibas.  
 780 Normalized object coordinate space for category-level 6d object pose and size estimation. In  
 781 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 2642–  
 782 2651, 2019.

783 Jianyuan Wang, Minghao Chen, Nikita Karaev, Andrea Vedaldi, Christian Rupprecht, and David  
 784 Novotny. Vggt: Visual geometry grounded transformer. In *Proceedings of the Computer Vision*  
 785 *and Pattern Recognition Conference*, pp. 5294–5306, 2025.

786 Shuzhe Wang, Vincent Leroy, Yohann Cabon, Boris Chidlovskii, and Jerome Revaud. Dust3r: Ge-  
 787 ometric 3d vision made easy. In *Proceedings of the IEEE/CVF Conference on Computer Vision*  
 788 *and Pattern Recognition*, pp. 20697–20709, 2024.

789 Ting-Chun Wang, Arun Mallya, and Ming-Yu Liu. One-shot free-view neural talking-head synthesis  
 790 for video conferencing. In *CVPR*, 2021.

791 Philippe Weinzaepfel, Vincent Leroy, Thomas Lucas, Romain Brégier, Yohann Cabon, Vaibhav  
 792 Arora, Leonid Antsfeld, Boris Chidlovskii, Gabriela Csurka, and Jérôme Revaud. Croco: Self-  
 793 supervised pre-training for 3d vision tasks by cross-view completion. *Advances in Neural Infor-  
 794 mation Processing Systems*, 35:3502–3516, 2022.

795 Yaniv Wolf, Amit Bracha, and Ron Kimmel. Gs2mesh: Surface reconstruction from gaussian splat-  
 796 ting via novel stereo views. In *European Conference on Computer Vision*, pp. 207–224. Springer,  
 797 2024.

798 Enze Xie, Wenhui Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Seg-  
 799 Former: Simple and efficient design for semantic segmentation with transformers. *Advances in*  
 800 *neural information processing systems*, 34:12077–12090, 2021.

801 Chao Xu, Ang Li, Linghao Chen, Yulin Liu, Ruoxi Shi, Hao Su, and Minghua Liu. Sparp: Fast  
 802 3d object reconstruction and pose estimation from sparse views. In *European Conference on*  
 803 *Computer Vision*, pp. 143–163. Springer, 2024.

804 Yuanwen Yue, Anurag Das, Francis Engelmann, Siyu Tang, and Jan Eric Lenssen. Improving 2D  
 805 Feature Representations by 3D-Aware Fine-Tuning. In *European Conference on Computer Vision*  
 806 (*ECCV*), 2024.

810 Junyi Zhang, Charles Herrmann, Junhwa Hur, Luisa Polania Cabrera, Varun Jampani, Deqing Sun,  
 811 and Ming-Hsuan Yang. A tale of two features: Stable diffusion complements dino for zero-shot  
 812 semantic correspondence. *Advances in Neural Information Processing Systems*, 36:45533–45547,  
 813 2023a.

814 Junyi Zhang, Charles Herrmann, Junhwa Hur, Eric Chen, Varun Jampani, Deqing Sun, and Ming-  
 815 Hsuan Yang. Telling left from right: Identifying geometry-aware semantic correspondence. In  
 816 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.  
 817 3076–3085, 2024a.

818 Junyi Zhang, Charles Herrmann, Junhwa Hur, Varun Jampani, Trevor Darrell, Forrester Cole, De-  
 819 qing Sun, and Ming-Hsuan Yang. Monst3r: A simple approach for estimating geometry in the  
 820 presence of motion. *arXiv preprint arXiv:2410.03825*, 2024b.

821 Longwen Zhang, Zijun Zhao, Xinzhou Cong, Qixuan Zhang, Shuqi Gu, Yuchong Gao, Rui Zheng,  
 822 Wei Yang, Lan Xu, and Jingyi Yu. Hack: Learning a parametric head and neck model for high-  
 823 fidelity animation. *ACM Transactions on Graphics (TOG)*, 42(4):1–20, 2023b.

824 Xiaozheng Zheng, Zhuo Su, Chao Wen, Zhou Xue, and Xiaojie Jin. Realistic full-body tracking  
 825 from sparse observations via joint-level modeling. In *Proceedings of the IEEE/CVF International  
 Conference on Computer Vision*, pp. 14678–14688, 2023a.

826 Yang Zheng, Adam W Harley, Bokui Shen, Gordon Wetzstein, and Leonidas J Guibas.  
 827 Pointodyssey: A large-scale synthetic dataset for long-term point tracking. In *Proceedings of  
 828 the IEEE/CVF International Conference on Computer Vision*, pp. 19855–19865, 2023b.

829 Yinglin Zheng, Hao Yang, Ting Zhang, Jianmin Bao, Dongdong Chen, Yangyu Huang, Lu Yuan,  
 830 Dong Chen, Ming Zeng, and Fang Wen. General facial representation learning in a visual-  
 831 linguistic manner. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern  
 832 Recognition*, pp. 18697–18709, 2022.

833 Hao Zhu, Wayne Wu, Wentao Zhu, Liming Jiang, Siwei Tang, Li Zhang, Ziwei Liu, and  
 834 Chen Change Loy. Celebvhq: A large-scale video facial attributes dataset. In *European con-  
 835 ference on computer vision*, pp. 650–667. Springer, 2022.

836 Junzhe Zhu, Yuanchen Ju, Junyi Zhang, Muhan Wang, Zhecheng Yuan, Kaizhe Hu, and Huazhe  
 837 Xu. Densematch: Learning 3d semantic correspondence for category-level manipulation from  
 838 a single demo. *arXiv preprint arXiv:2412.05268*, 2024.

839 Xiangyu Zhu, Xiaoming Liu, Zhen Lei, and Stan Z Li. Face alignment in full pose range: A 3d total  
 840 solution. *IEEE transactions on pattern analysis and machine intelligence*, 41(1):78–92, 2017.

841 Wojciech Zienonka, Timo Bolkart, and Justus Thies. Towards metrical reconstruction of human  
 842 faces. In *European conference on computer vision*, pp. 250–269. Springer, 2022.

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863