

000 STREAMING AUTOREGRESSIVE VIDEO GENERA- 001 002 003 004 005 006 007 008 009 010 011 012 013 014 015 016 017 018 019 020 021 022 023 024 025 026 027 028 029 030 031 032 033 034 035 036 037 038 039 040 041 042 043 044 045 046 047 048 049 050 051 052 053 054 055 056 057 058 059 060 061 062 063 064 065 066 067 068 069 070 071 072 073 074 075 076 077 078 079 080 081 082 083 084 085 086 087 088 089 090 091 092 093 094 095 096 097 098 099 100 101 102 103 104 105 106 107 108 109 110 111 112 113 114 115 116 117 118 119 120 121 122 123 124 125 126 127 128 129 130 131 132 133 134 135 136 137 138 139 140 141 142 143 144 145 146 147 148 149 150 151 152 153 154 155 156 157 158 159 160 161 162 163 164 165 166 167 168 169 170 171 172 173 174 175 176 177 178 179 180 181 182 183 184 185 186 187 188 189 190 191 192 193 194 195 196 197 198 199 200 201 202 203 204 205 206 207 208 209 210 211 212 213 214 215 216 217 218 219 220 221 222 223 224 225 226 227 228 229 230 231 232 233 234 235 236 237 238 239 240 241 242 243 244 245 246 247 248 249 250 251 252 253 254 255 256 257 258 259 260 261 262 263 264 265 266 267 268 269 270 271 272 273 274 275 276 277 278 279 280 281 282 283 284 285 286 287 288 289 290 291 292 293 294 295 296 297 298 299 300 301 302 303 304 305 306 307 308 309 310 311 312 313 314 315 316 317 318 319 320 321 322 323 324 325 326 327 328 329 330 331 332 333 334 335 336 337 338 339 340 341 342 343 344 345 346 347 348 349 350 351 352 353 354 355 356 357 358 359 360 361 362 363 364 365 366 367 368 369 370 371 372 373 374 375 376 377 378 379 380 381 382 383 384 385 386 387 388 389 390 391 392 393 394 395 396 397 398 399 400 401 402 403 404 405 406 407 408 409 410 411 412 413 414 415 416 417 418 419 420 421 422 423 424 425 426 427 428 429 430 431 432 433 434 435 436 437 438 439 440 441 442 443 444 445 446 447 448 449 450 451 452 453 454 455 456 457 458 459 460 461 462 463 464 465 466 467 468 469 470 471 472 473 474 475 476 477 478 479 480 481 482 483 484 485 486 487 488 489 490 491 492 493 494 495 496 497 498 499 500 501 502 503 504 505 506 507 508 509 510 511 512 513 514 515 516 517 518 519 520 521 522 523 524 525 526 527 528 529 530 531 532 533 534 535 536 537 538 539 540 541 542 543 544 545 546 547 548 549 550 551 552 553 554 555 556 557 558 559 559 560 561 562 563 564 565 566 567 568 569 569 570 571 572 573 574 575 576 577 578 579 579 580 581 582 583 584 585 586 587 588 589 589 590 591 592 593 594 595 596 597 598 599 599 600 601 602 603 604 605 606 607 608 609 609 610 611 612 613 614 615 616 617 618 619 619 620 621 622 623 624 625 626 627 628 629 629 630 631 632 633 634 635 636 637 638 639 639 640 641 642 643 644 645 646 647 648 649 649 650 651 652 653 654 655 656 657 658 659 659 660 661 662 663 664 665 666 667 668 669 669 670 671 672 673 674 675 676 677 678 679 679 680 681 682 683 684 685 686 687 688 689 689 690 691 692 693 694 695 696 697 698 699 699 700 701 702 703 704 705 706 707 708 709 709 710 711 712 713 714 715 716 717 718 719 719 720 721 722 723 724 725 726 727 728 729 729 730 731 732 733 734 735 736 737 738 739 739 740 741 742 743 744 745 746 747 748 749 749 750 751 752 753 754 755 756 757 758 759 759 760 761 762 763 764 765 766 767 768 769 769 770 771 772 773 774 775 776 777 778 779 779 780 781 782 783 784 785 786 787 788 789 789 790 791 792 793 794 795 796 797 798 799 799 800 801 802 803 804 805 806 807 808 809 809 810 811 812 813 814 815 816 817 818 819 819 820 821 822 823 824 825 826 827 828 829 829 830 831 832 833 834 835 836 837 838 839 839 840 841 842 843 844 845 846 847 848 849 849 850 851 852 853 854 855 856 857 858 859 859 860 861 862 863 864 865 866 867 868 869 869 870 871 872 873 874 875 876 877 878 879 879 880 881 882 883 884 885 886 887 888 889 889 890 891 892 893 894 895 896 897 898 899 899 900 901 902 903 904 905 906 907 908 909 909 910 911 912 913 914 915 916 917 918 919 919 920 921 922 923 924 925 926 927 928 929 929 930 931 932 933 934 935 936 937 938 939 939 940 941 942 943 944 945 946 947 948 949 949 950 951 952 953 954 955 956 957 958 959 959 960 961 962 963 964 965 966 967 968 969 969 970 971 972 973 974 975 976 977 978 979 979 980 981 982 983 984 985 986 987 988 989 989 990 991 992 993 994 995 996 997 998 999 1000 1001 1002 1003 1004 1005 1006 1007 1008 1009 1009 1010 1011 1012 1013 1014 1015 1016 1017 1018 1019 1019 1020 1021 1022 1023 1024 1025 1026 1027 1028 1029 1029 1030 1031 1032 1033 1034 1035 1036 1037 1038 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 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 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 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 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 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 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 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 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 1898 1899 1899 1900 1901 1902 1903 1904 1905 1906 1907 1908 1909 1909 1910



Figure 1: Our Diagonal Distillation framework achieves comparable quality to the full-step model while significantly reducing latency. The method yields a 1.88 \times speedup on 5-second short video generation on a single H100 GPU.

into efficient few-step sampling AR model. Recent training methods (Chen et al., 2024; Gao et al., 2024b; Gu et al., 2025; Hu et al., 2024; Li et al., 2024b; Liu et al., 2024b; Weng et al., 2024; Yin et al., 2025; Zhang et al., 2025a;b) have further improved stability and efficiency, making interactive applications increasingly feasible (Arriola et al., 2025; Liu et al., 2024c).

Despite these advances, existing video distillation methods are largely adapted from image generation, and their direct extension to video often yields suboptimal results. This limitation arises from insufficient consideration of the temporal dimension and the neglect of inter-frame consistency. As a result, multi-step sampling remains essential for maintaining high-quality video generation. For example, while autoregressive frameworks such as Causvid (Yin et al., 2025) and Self Forcing (Huang et al., 2025) can reduce latency, they still require multiple steps per segment, and compressing them to fewer steps leads to noticeable performance degradation.

Our guiding insight is that, in autoregressive video generation, predicting the next chunk inherently requires predicting the next noise level (see Figure 2). This implicit prediction, however, introduces two critical challenges. First, autoregressive video models often suffer from exposure bias. When predicting the next chunk conditioned on previously generated clean frames, the model must implicitly predict the next noise level for subsequent frames. This can lead to progressive degradation, such as over-saturation in later frames, as errors in noise-level prediction accumulate over time. Although techniques like Self Forcing (Huang et al., 2025) have been proposed to mitigate exposure bias by using model-generated content during training, they still struggle to maintain visual quality over long sequences. Second, the same phenomenon implies that if structural priors are captured in early chunks, later chunks can generate relatively clear frames even with fewer denoising steps. However, existing distillation approaches often discard valuable temporal context accumulated across denoising steps in video generation models, which is essential for preserving coherence and detail when reducing the sampling steps.

Motivated by these insights, we introduce a flow-aware diagonal distillation framework – **DiaDistill** that redefines the temporal context incorporation by leveraging information across both time and denoising steps. Departing from standard practices that process chunks in isolation, our method employs a novel diagonal forcing mechanism operating jointly across time and denoising steps. This results in a diagonal denoising trajectory wherein earlier chunks are denoised with more steps, while later chunks use progressively fewer. This strategy improves computational efficiency by using less denoising steps in total and allows each chunk to inherit denoising trajectories from prior chunks as contextual priors—a training paradigm we term Diagonal Forcing. By explicitly simulating diagonal denoising paths during training through controlled noise injection, Diagonal Forcing enhances self-conditioned generation and mitigates error accumulation in long videos. Furthermore, we empirically observe that employing very few steps in later chunks can attenuate motion amplitude. To counteract this, we introduce Flow Distribution Matching, which integrates explicit temporal modeling into the distillation loss. This approach preserves dynamic consistency by ensuring the predicted motion distributions align with those of the full-step model, thus ensuring that the student model not only matches the teacher in image quality but also faithfully preserves motion characteristics. The contributions of this work are:

- We propose **Diagonal Distillation**, a method for high-quality video generation during model distillation and inference. It allocates more denoising steps to earlier chunks and progressively fewer to later ones, rather than keeping the number of steps constant across all chunks. This



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

Figure 2: We find that when the training data uses explicit noise frames as conditions in Causvid (Yin et al., 2025), the next chunk prediction essentially functions as an implicit next noise level prediction. It can be observed that even with single-step prediction, the image progressively becomes clearer.

- approach achieves an improved trade-off between quality and efficiency by leveraging contextual structured priors in AR video generation.
- We introduce **Diagonal Forcing** based on diagonal distillation, a unified method that operates along both temporal and denoising-step dimensions. It leverages trajectories from prior chunks as contextual priors and explicitly simulates diagonal denoising paths during training through controlled noise injection, thereby reducing long-term error accumulation.
 - We present **Flow Distribution Matching** as a complementary component to diagonal distillation, designed to mitigate motion degradation and amplitude attenuation in later chunks. By incorporating explicit temporal modeling into the distillation loss, this approach enhances dynamic consistency and ensures smooth motion transitions.

2 RELATED WORK

Diffusion Distillation. Diffusion distillation accelerates sampling via deterministic or distributional approaches. Deterministic methods (e.g., progressive distillation (Salimans & Ho, 2022), consistency distillation (Li et al., 2023; Song et al., 2023), rectified flow (Lamb et al., 2016)) regress noise-to-sample mappings but often yield blurry outputs with few steps due to optimization challenges (Kingma et al., 2021), typically requiring multiple steps (e.g., eight) for acceptable quality (Li et al., 2023; 2024a). Distributional methods approximate the teacher’s distribution using adversarial training (Brooks et al., 2024; Ho et al., 2022), score distillation (Li et al., 2022; Luo et al., 2024), or hybrid objectives. Recent hybrids combine both paradigms but still suffer from one-step artifacts and commonly need multi-step sampling. Representative works include LADD (Sauer et al., 2024a) (relies on expensive pre-generated teacher targets), Lightning (Lin et al., 2024) and Hyper (Ren et al., 2024) (require intermediate timestep supervision), and DMD/DMD2 (Yin et al., 2024b;a) and ADD (Sauer et al., 2024b) (integrate adversarial and score matching losses). While these distillation methods have shown impressive results in image generation, their direct application to video often yields suboptimal results due to insufficient consideration of the temporal dimension and inter-frame consistency. Our work addresses this gap by proposing a flow-aware diagonal distillation framework specifically designed for video generation, which leverages temporal context across both time and denoising steps to maintain coherence while reducing sampling steps.

Autoregressive, Diffusion, and Hybrid Video Generation. Modern video generation is dominated by scalable diffusion and autoregressive (AR) models. Video diffusion models use bidirectional attention to denoise all frames concurrently (Blattmann et al., 2023a;b; Brooks et al., 2024; Deng et al., 2024; Kong et al., 2024; Polyak et al., 2024; Villegas et al., 2022; Wan et al., 2025; Yang et al., 2024), while AR models generate spatiotemporal tokens sequentially via next-token prediction (Bruce et al., 2024; Kondratyuk et al., 2023; Ren et al., 2025; Wang et al., 2024; Weissenborn et al., 2019; Yan et al., 2021; Liu et al., 2025). Hybrid models that merge these two paradigms have recently emerged as a promising direction (Chen et al., 2024; Gao et al., 2024b; Gu et al., 2025; Hu et al., 2024; Jin et al., 2024; Li et al., 2024b; Liu et al., 2024a;b; Weng et al., 2024; Yin et al., 2025; Zhang et al., 2025a;b), also in other sequence domains (Arriola et al., 2025; Liu et al., 2024c). These hybrids typically integrate diffusion into AR generation to boost visual quality, but they still require multiple denoising steps per chunk, hindering real-time deployment. Our work builds on these hybrids, drawing inspiration from Yin et al. (2025) and Huang et al. (2025) to mitigate exposure bias. However, these methods still face challenges with long-term error accumulation and motion degradation when compressed to fewer steps. Our proposed *Diagonal Distillation* framework addresses these issues via a novel diagonal forcing mechanism operating jointly across time and denoising steps, enabling efficient computation while preserving temporal coherence. The *Diagonal*

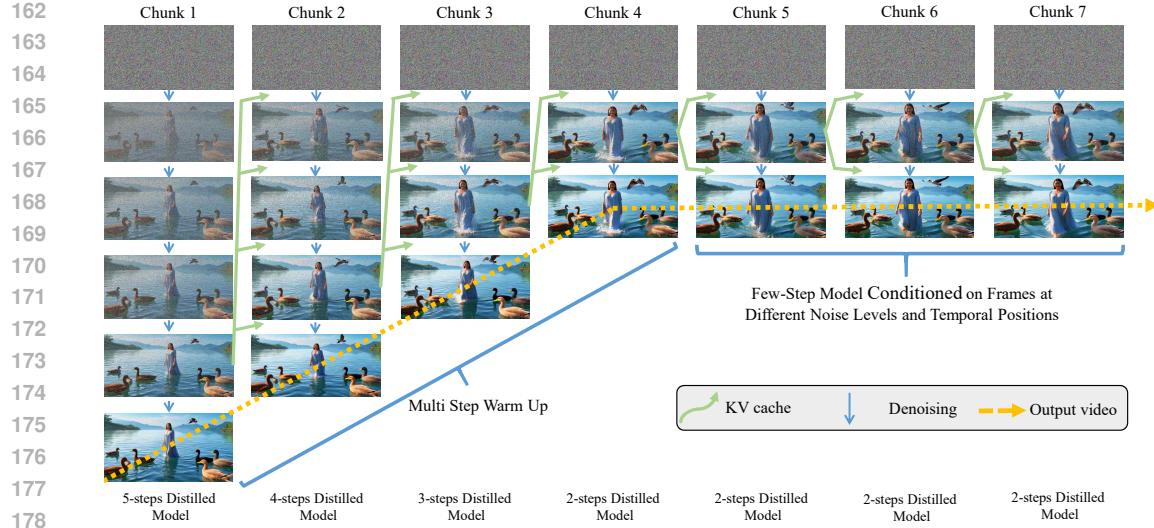


Figure 3: Diagonal Denoising with Diagonal Forcing and Progressive Step Reduction. We illustrate our method starting with 5 denoising steps for the first chunk and progressively reducing them to 2 steps by Chunk 7 (see Section 4.3 for more parameters). For chunks with $k \geq 4$, we use a fixed two-step denoising process, reusing the Key-Value (KV) cache from the previous chunk's last noisy frame. This approach maintains temporal coherence while reducing latency, the pseudo-code is provided in the appendix.

Forcing training paradigm explicitly simulates diagonal denoising paths to enhance self-conditioned generation, and *Flow Distribution Matching* ensures motion consistency with reduced steps.

3 METHODOLOGY

3.1 PRELIMINARY AND FRAMEWORK OVERVIEW

Diffusion Models generate data through an iterative denoising process. The forward diffusion process progressively corrupts a sample $x \sim p_{\text{real}}$ over T steps, such that at timestep t , the diffused sample follows $p_{\text{real},t}(x_t) = \int p_{\text{real}}(x)q(x_t|x)dx$, with $q_t(x_t|x) \sim \mathcal{N}(\alpha_t x, \sigma_t^2 I)$, where $\alpha_t, \sigma_t > 0$ are determined by the noise schedule. The model learns to reverse this process by predicting a denoised estimate $\mu(x_t, t)$. The score function of the diffused distribution is:

$$s_{\text{real}}(x_t, t) = \nabla_{x_t} \log p_{\text{real},t}(x_t) = -\frac{x_t - \alpha_t \mu_{\text{real}}(x_t, t)}{\sigma_t^2}. \quad (1)$$

Sampling typically requires many iterative steps. Distribution Matching Distillation (DMD) distills a multi-step diffusion model (teacher) into a one-step generator G by minimizing the KL divergence between the diffused real and generated distributions, $p_{\text{real},t}$ and $p_{\text{fake},t}$. The gradient of this loss is:

$$\nabla \mathcal{L}_{\text{DMD}} = \mathbb{E}_t (\nabla_{\theta} \text{KL}(p_{\text{fake},t} \| p_{\text{real},t})) = -\mathbb{E}_t \left(\int (s_{\text{real}}(F(G_{\theta}(z), t), t) - s_{\text{fake}}(F(G_{\theta}(z), t), t)) \frac{dG_{\theta}(z)}{d\theta} dz \right), \quad (2)$$

where $z \sim \mathcal{N}(0, \mathbf{I})$, F is the forward diffusion process, and $s_{\text{real}}, s_{\text{fake}}$ are scores from models trained on real and generated data. An additional regression loss is often used for regularization:

$$\mathcal{L}_{\text{reg}} = E_{(z,y)} d(G_{\theta}(z), y), \quad (3)$$

where y is an image generated by the teacher from z . Directly applying DMD to video generation faces a significant challenge: the regression loss \mathcal{L}_{reg} primarily ensures per-frame quality but fails to explicitly capture the underlying temporal coherence and long-range dependencies between frames, which are critical for video quality. This often results in degraded fluidity and consistency. To overcome this, we extend the DMD framework with two core innovations: 1) a **Diagonal Denoising with Diagonal Forcing** strategy that manages long-sequence generation and reduces error accumulation (Section 3.2) a novel **Flow Distribution Matching** objective that explicitly aligns the temporal dynamics of the student and teacher models (Section 3.3).

216 3.2 DIAGONAL DENOISING WITH DIAGONAL FORCING
 217

218 Building upon the DMD foundation, we present diagonal distillation, a framework for efficient video
 219 generation. As illustrated in Figure 3, our approach introduces a Diagonal Denoising strategy that
 220 progressively reduces denoising steps across video chunks, combined with a novel Diagonal Forcing
 221 mechanism to maintain temporal coherence and mitigate error accumulation.

222 **Diagonal Denoising: Progressive Step Reduction Strategy** Our core innovation is a diagonal
 223 denoising strategy that allocates computation based on temporal importance. The method assigns
 224 more denoising steps to earlier chunks and progressively fewer to later ones, rather than maintaining
 225 a constant number of steps across all chunks. This approach achieves an improved trade-off between
 226 quality and efficiency by leveraging contextual structured priors in autoregressive video generation.
 227 For the first three chunks ($k = 1, 2, 3$), we use distilled models with decreasing steps ($s_k = 5, 4, 3$):

$$\mathbf{X}_k = \mathcal{D}_{s_k}(\mathbf{Z}_k | \tilde{\mathbf{X}}_{<k}), \quad (4)$$

230 where \mathbf{X}_k is the k -th chunk output, $\mathbf{Z}_k \sim \mathcal{N}(0, \mathbf{I})$ is Gaussian noise, and $\tilde{\mathbf{X}}_{<k}$ contains previously
 231 noised chunks. For $k \geq 4$, we employ efficient two-step denoising:

$$\mathbf{C}_k = \mathcal{T}(\tilde{\mathbf{X}}_{k-1}), \mathbf{X}_k = \mathcal{D}_2(\mathcal{D}_1(\mathbf{Z}_k | \mathbf{C}_k) | \mathbf{C}_k), \quad (5)$$

232 where \mathbf{C}_k is the conditioning signal derived from previous chunks, \mathcal{T} denotes the conditioning
 233 module, and $\mathcal{D}_1, \mathcal{D}_2$ represent the first and second denoising steps respectively.

234 **Diagonal Forcing: Contextual Prior Propagation** The core innovation of Diagonal Forcing
 235 lies in its explicit modeling of diagonal denoising trajectories during training through controlled
 236 noise injection. This approach ensures temporal coherence across chunks while minimizing error
 237 accumulation by conditioning each new chunk on the final noised state from the previous chunk's
 238 diffusion process. Specifically, the conditioning input for chunk k is derived from the clean output
 239 \mathbf{X}_{k-1} of chunk $k - 1$ through a noise injection operation:

$$\tilde{\mathbf{X}}_{k-1} = \sqrt{\alpha_{k-1}} \mathbf{X}_{k-1} + \sqrt{1 - \alpha_{k-1}} \cdot \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (6)$$

240 where α_{k-1} controls the noise schedule along the diagonal path and $\boldsymbol{\epsilon}$ is standard Gaussian noise.
 241 This formulation explicitly maintains the diagonal denoising trajectory $\mathbf{X}_k \rightarrow \tilde{\mathbf{X}}_{k-1} \rightarrow \mathbf{X}_{k-1}$,
 242 where $\tilde{\mathbf{X}}_{k-1}$ serves as the KV cache input for chunk k . By propagating these noised representations
 243 across chunks, the method effectively leverages denoising trajectories from prior chunks as contextual
 244 priors. The diagonal alignment of these trajectories ensures that error accumulation is minimized
 245 while preserving long-range coherence in the generated output.

246 3.3 FLOW DISTRIBUTION MATCHING
 247

248 Motion attenuation in few-step denoising stems from truncated noise estimation paths. We quantify
 249 the temporal distribution mismatch through flow-based divergence:

$$\mathcal{E}_{\text{motion}} = D_{\text{KL}}(p_{\text{teacher}}(\mathcal{F}(\mathbf{x}) | \mathbf{x}_t) \| p_{\text{student}}(\mathcal{F}(\mathbf{x}) | \mathbf{x}_t)) \quad (7)$$

250 where $\mathcal{F}(\mathbf{x})$ represents the motion flow field extracted from video sequence \mathbf{x} . This measures the
 251 distributional divergence between teacher and student in the temporal dimension.

252 The standard Distribution Matching Distillation (DMD) framework minimizes spatial divergence
 253 through reverse KL minimization. We extend this to the temporal domain by defining flow distribution
 254 matching:

$$\nabla_{\phi} \mathcal{L}_{\text{DMD}}^{\text{flow}} \triangleq \mathbb{E}_t (\nabla_{\phi} \text{KL}(p_{\text{gen,flow},t} \| p_{\text{data,flow},t})) \quad (8)$$

255 where $p_{\text{data,flow},t} = p(\mathcal{F}(\mathbf{x}) | \Psi(\mathbf{x}, t))$ is the smoothed flow distribution from real data, and
 256 $p_{\text{gen,flow},t} = p(\mathcal{F}(\mathbf{x}) | \Psi(G_{\phi}(\boldsymbol{\epsilon}), t))$ is the generator's flow distribution. The gradient approximation
 257 for flow distribution matching follows the DMD framework:

$$\nabla_{\phi} \mathcal{L}_{\text{DMD}}^{\text{flow}} \approx -\mathbb{E}_t \left[\int (s_{\text{data}}^{\text{flow}}(\Psi(G_{\phi}(\boldsymbol{\epsilon}), t), t) - s_{\text{gen},\xi}^{\text{flow}}(\Psi(G_{\phi}(\boldsymbol{\epsilon}), t), t)) \frac{dG_{\phi}(\boldsymbol{\epsilon})}{d\phi} d\boldsymbol{\epsilon} \right], \quad (9)$$

258 where $s_{\text{data}}^{\text{flow}}$ and $s_{\text{gen},\xi}^{\text{flow}}$ are the flow score functions defined as:
 259

$$s^{\text{flow}}(\mathbf{x}_t, t) = \nabla_{\mathbf{x}_t} \log p(\mathcal{F}(\mathbf{x}) | \mathbf{x}_t). \quad (10)$$

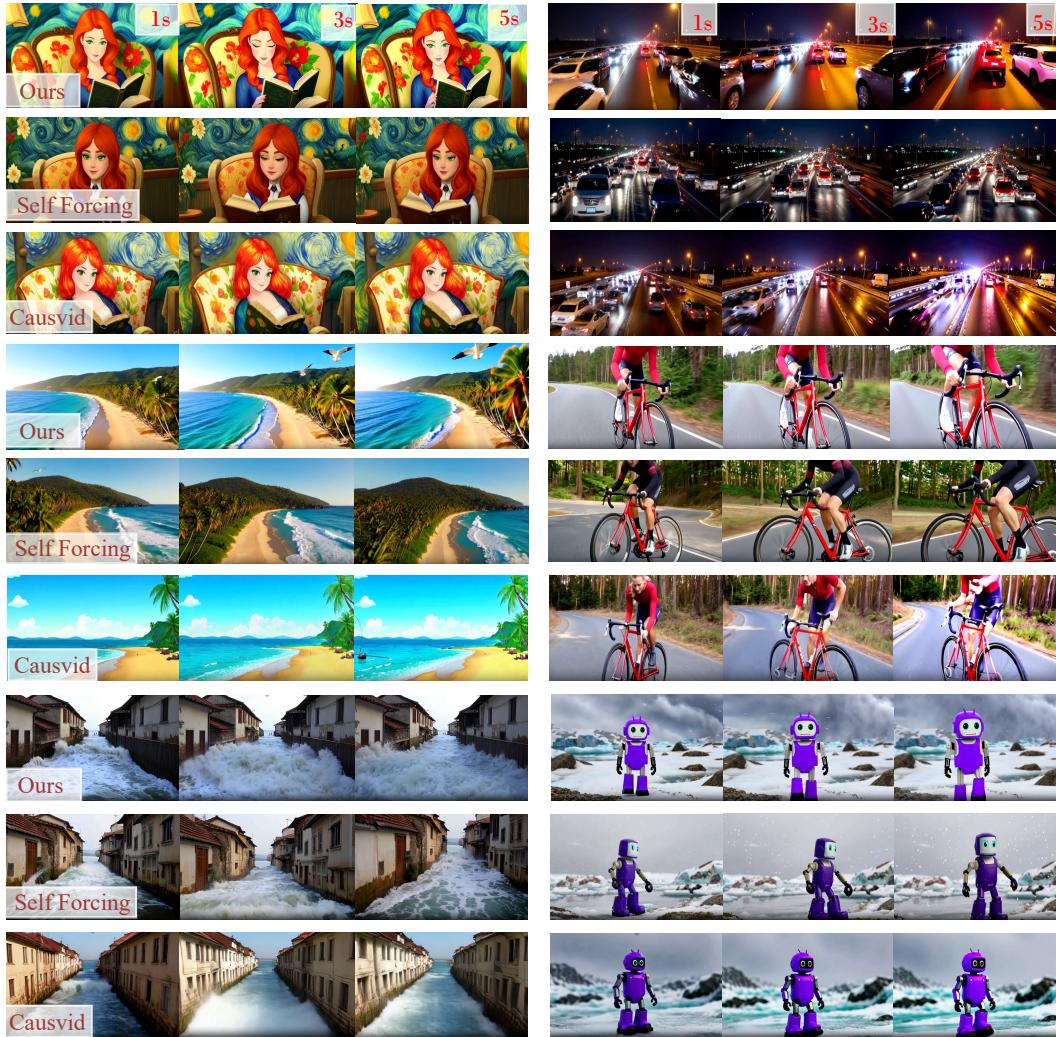


Figure 4: Comparing the results from three different models. For more results, please refer to our supplementary material.

To operationalize this framework, we employ a flow regression loss for feature alignment:

$$\mathcal{L}_{\text{reg}}^{\text{flow}} = \mathbb{E}_{t,\epsilon} [\|\mathcal{F}(G_{\phi}^{\text{teacher}}(\epsilon, t)) - \mathcal{F}(G_{\phi}^{\text{student}}(\epsilon, t))\|_2^2], \quad (11)$$

where $G_{\phi}(\epsilon, t)$ denotes the generator output at timestep t . Our method uses a lightweight, self-contained motion feature extraction module $\mathcal{F}(\cdot)$ that operates directly on latent representations, avoiding dependencies on external pre-trained optical flow estimators. Specifically, we implement $\mathcal{F}(\cdot)$ as a learnable representation with convolution on latent difference: it first computes the difference between consecutive latent frames, then applies convolutional layers to extract local motion patterns, followed by an MLP for feature adaptation. The student version is trainable with gradient flow, while the teacher components are updated via EMA, ensuring stable and efficient motion representation learning. The overall training objective combines both spatial and temporal distribution matching:

$$\mathcal{L}_{\text{Total}} = \lambda_{\text{spatial}} \mathcal{L}_{\text{DMD}} + \mathcal{L}_{\text{reg}} + \gamma (\lambda_{\text{flow}} \mathcal{L}_{\text{DMD}}^{\text{flow}} + \mathcal{L}_{\text{reg}}^{\text{flow}}), \quad (12)$$

where γ weights the temporal terms. Where we set $\lambda_{\text{spatial}}=4$ and $\lambda_{\text{flow}}=4$ This framework jointly minimizes motion distribution divergence while maintaining spatial fidelity in the distilled video model.

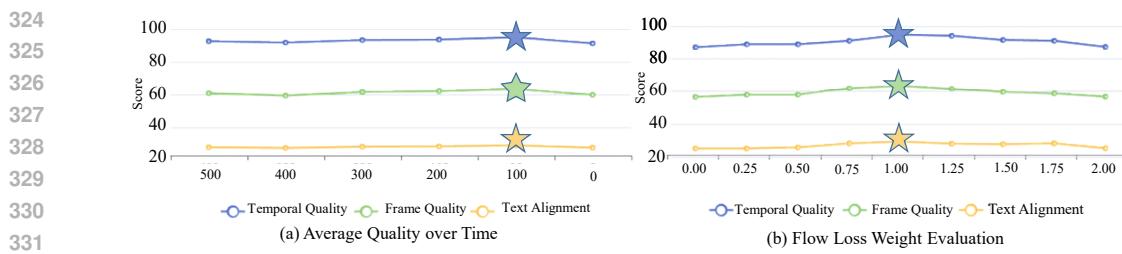


Figure 5: Ablation study results:(a)Performance evaluation across different diagonal forcing timesteps, demonstrating optimal outcomes at 100 steps (1000 steps correspond to complete noise addition, while 0 steps represent the clean frame);(b)Impact of motion loss weight on model performance.

4 EXPERIMENTS AND RESULTS

4.1 IMPLEMENTATION DETAILS

Training Details We implemented DiaDistill using Wan2.1-T2V-1.3B (Wan et al., 2025), a model based on Flow Matching (Lipman et al., 2022) that is capable of generating 5 videos at 16 FPS with a resolution of 832×480 . For both ODE initialization and Diagonal Distillation training, we sampled text prompts from a filtered and LLM-extended version of VidProM (Wang & Yang, 2024).

Inference Details To assess real-time applicability, we measured both throughput (frames per second) and first-frame latency, acknowledging that true real-time performance requires exceeding video playback rates while maintaining imperceptible delay. All speed tests were conducted on a single NVIDIA H100 GPU with tiny VAE Boer Bohan (2025). The core component is the rolling KV cache mechanism following Self-Forcing Huang et al. (2025), which operates with a chunk size of 3 frames. Our buffering strategy is implemented using a fixed-size KV cache that maintains context from the most recent 4 chunks, resulting in a consistent memory footprint of 17.5 GB. For detailed ablation analysis please refer to our supplementary.

Evaluation Details We evaluated visual quality and semantic consistency using VBench (Huang et al., 2024). Temporal Quality is the average of Subject Consistency, Background Consistency, Temporal Flickering, Motion Smoothness, and Dynamic Degree. Frame Quality is the average of Aesthetic Quality and Imaging Quality. Text Alignment is the average of Object Class, Multiple Objects, Human Action, Color, Spatial Relationship, Scene, Appearance, Style, and Temporal Style. The aggregation method for each score is a simple arithmetic mean of the normalized scores from its constituent sub-dimensions. This approach is consistent with prior works like Causvid and Self-Forcing for fair comparison.

4.2 COMPARISON WITH STATE-OF-THE-ART METHODS

We evaluate DiaDistill against five state-of-the-art video generation methods: Wan2.1 (Wan et al., 2025), SkyReels-V2 (Chen et al., 2025), MAGI-1 (Teng et al., 2025), Causvid (Yin et al., 2025), and Self Forcing (Huang et al., 2025).

| Model | Throughput↑ | First-Frame Latency ↓ | Speedup | Total↑ | Quality↑ | Semantic↑ |
|-----------------------------------|-------------|-----------------------|---------------|--------------|-------------|--------------|
| Wan2.1 (Wan et al., 2025) | 0.78 | 103 | 1.0× | 84.26 | 85.3 | 80.09 |
| SkyReels-V2 (Chen et al., 2025) | 0.49 | 112 | 0.91× | 82.67 | 84.70 | 74.53 |
| MAGI-1 (Teng et al., 2025) | 0.19 | 282 | 0.36× | 79.18 | 82.04 | 67.74 |
| Causvid (Yin et al., 2025) | 17.0 | 0.69 | 149.3× | 81.20 | 84.05 | 69.80 |
| Self Forcing (Huang et al., 2025) | 17.0 | 0.69 | 149.3× | 84.31 | 85.07 | 81.28 |
| DiaDistill (Ours) | 31.0 | 0.37 | 277.3× | 84.48 | 85.26 | 81.73 |

Table 1: Comprehensive comparison of video generation methods

Forcing (Huang et al., 2025). As shown in Table 1, our method achieves a $277.3\times$ speedup over the Wan2.1 baseline while maintaining competitive visual quality (85.26 vs. 85.3). This represents a $1.53\times$ improvement in latency over the previous fastest method, Self Forcing ($149.3\times$), alongside superior overall performance and semantic consistency . Qualitative results in Figure 4 further demonstrate advantages in temporal consistency, with smoother frame transitions and fewer dynamic artifacts. Visual fidelity improvements are most apparent in complex motions and textures, where baseline methods exhibit blurring or distortion. These findings collectively show that DiaDistill effectively balances the traditional trade-off between generation quality and computational efficiency.

4.3 ABLATION STUDIES

Key Components Diagonal Denoising assigns more denoising steps to early video chunks to establish a high-quality foundation and progressively reduces denoising steps for subsequent chunks, whereas

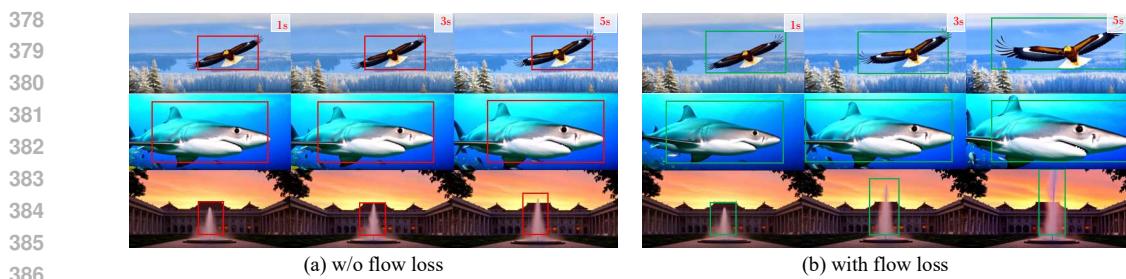


Figure 6: **Visual comparison of motion effects:** (a) Without motion loss shows minimal motion amplitude with only slight object movement; (b) With motion loss demonstrates significantly increased motion amplitude throughout the entire frame, validating our method’s effectiveness.

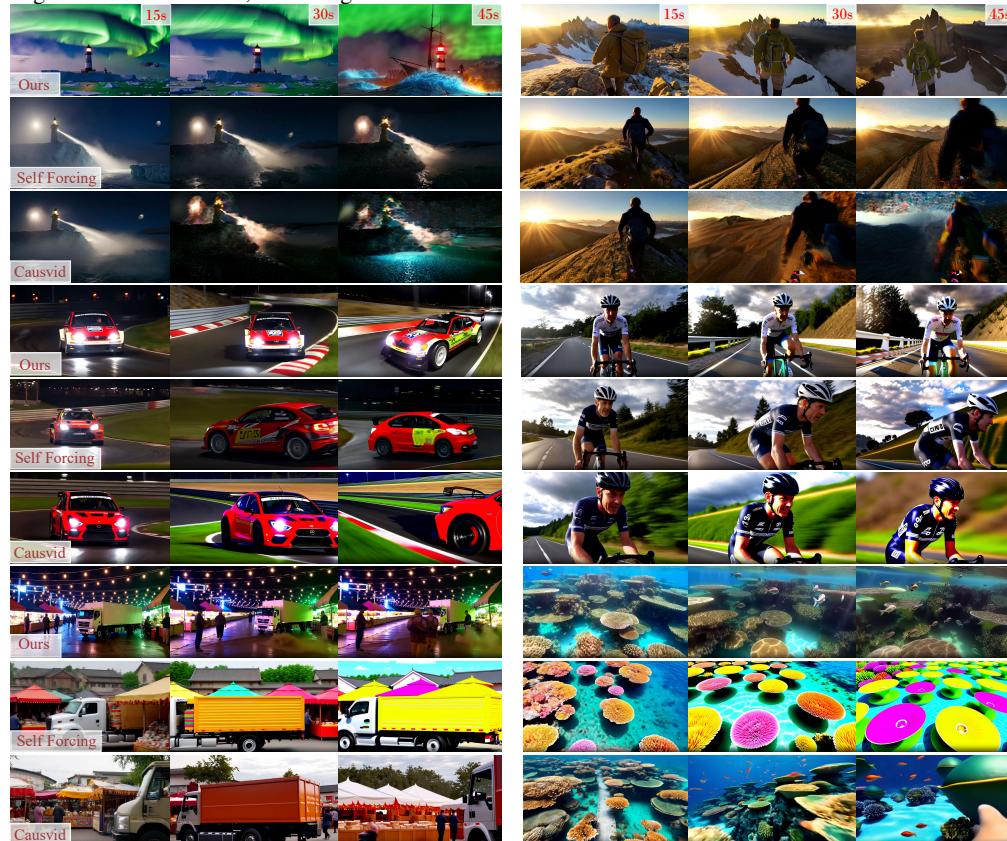


Figure 7: Qualitative comparison of long video generation(45s) with Self Forcing and Causvid. The visual results show that other methods suffer from noticeable saturation distortion and quality decay over time, whereas our approach preserves detail and consistency. Additional results are provided in the supplementary material.

without it, the same number of steps is applied uniformly across all chunks. Diagonal Forcing refers to using noisy frames instead of clean frames as the Key-Value (KV) cache in autoregressive generation. Our ablation study shows that removing either flow distribution matching loss or Diagonal Forcing significantly degrades video quality across all metrics (Table 2). Without Diagonal Denoising—which corresponds to the inference cost of Self Forcing in Table 1—we observe that the model achieves performance comparable to ours, but our method achieves a $1.53\times$ speedup. Notably, we find that flow distribution matching loss primarily benefits the few-step denoising regime and helps align its performance with the many-step denoising baseline (i.e., without Diagonal Denoising), and provides limited benefits when applied to a many-step denoising setting.

Diagonal Forcing Timesteps Moreover, we systematically evaluated diagonal forcing using metrics across different noise levels of timesteps for the kv cache. As Figure 5(a) shows, 100 timesteps achieved optimal scores across all evaluation dimensions, including temporal quality, frame quality, and text alignment. The performance peaks at this specific noise level before degrading as timesteps approach complete noise addition (1000 steps) or clean frames (0 steps). This can be attributed to the

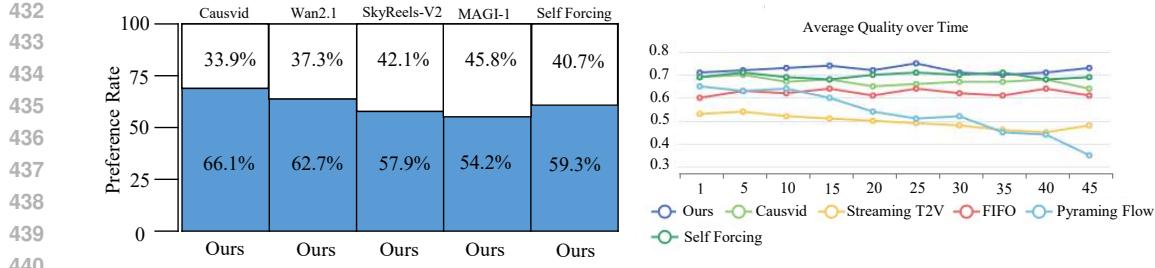


Figure 8: Quantitative evaluation of long video generation. The plot compares human preference scores and quality consistency over time for different methods under identical conditions. Our approach maintains stable quality throughout extended sequences, achieving scores above 50%, and attains a significant reduction in inference latency.

fact that excessive noise (high timesteps) blurs the structural priors in the video context, leading to reduced motion magnitude. This also explains why our method generates larger motion amplitudes compared to MAGI (Teng et al., 2025). Conversely, insufficient noise (low timesteps) causes the next chunk prediction to implicitly perform next noise level prediction, which can result in over-denoising of subsequent chunks and ultimately lead to over-saturated outputs.

| Ablation Variant | Temporal Quality ↑ | Frame Quality ↑ | Text Alignment ↑ | Total Score ↑ | Steps | Temporal Quality ↑ | Frame Quality ↑ | Text Alignment ↑ | NFEs | In-Flight Latency (s) ↓ | Throughput (FPS) ↑ |
|----------------------------|--------------------|-----------------|------------------|---------------|---------|--------------------|-----------------|------------------|------|-------------------------|--------------------|
| Without Diagonal Forcing | 92.1 | 60.1 | 26.9 | 83.58 | 4322222 | 94.9 | 63.4 | 28.9 | 34 | 0.23 ± 0.02 | 31.0 |
| Without Flow Loss | 92.5 | 60.8 | 27.8 | 84.18 | 5433333 | 95.1 | 63.2 | 29.3 | 48 | 0.34 ± 0.02 | 23.3 |
| Without Diagonal Denoising | 95.1 | 63.2 | 28.6 | 84.46 | 5432222 | 94.8 | 63.1 | 29.0 | 40 | 0.23 ± 0.02 | 29.7 |
| Full Method (Ours) | 94.9 | 63.4 | 28.9 | 84.48 | 5333333 | 95.0 | 63.9 | 29.1 | 46 | 0.34 ± 0.02 | 22.5 |
| | | | | | 4222222 | 95.0 | 63.7 | 28.5 | 44 | 0.34 ± 0.02 | 23.5 |
| | | | | | | 93.4 | 62.3 | 27.8 | 32 | 0.23 ± 0.02 | 32.0 |

Table 2: Ablation Study on Key Components of DiaDistill.

Table 3: Performance evaluation of denoising step configurations

Flow Loss Weight We conducted a comprehensive ablation study across eight motion loss weight configurations. Figure 5(b) reveals the crucial balance between motion guidance (via Flow Distribution Matching) and the DMD learning objectives, with optimal performance observed at a weight of 1.0. This balanced weighting scheme ensures the harmonious optimization of temporal consistency, frame quality, and textual alignment metrics.

Denoising Configurations We evaluated six denoising configurations (represented by 7-digit sequences specifying steps per chunk as a 5 seconds video have 7 chunks in our setting) across quality and computational metrics. As shown in Table 3, these configurations exhibit trade-offs between generation quality and efficiency. Among them, configuration 5333333 achieves the highest quality, while 4222222 offers the maximum throughput. To balance video quality and real-time performance, we selected configuration 4322222, as it has the second-lowest number of NFEs and delivers performance comparable to configurations with significantly higher latency and throughput, with only marginal differences.

KV Cache Scaling Analysis We further analyze the trade-offs between attention window size, video quality (Total Score), and in-flight interaction latency—the delay in responding to a new input signal during steady-state generation. As shown in Table 4, a larger window size generally improves quality at the cost of higher memory usage and, crucially, a longer delay in responding to user interactions. Performance plateaus around a window size of 18–27, leading us to select an optimal size that balances responsiveness with quality and efficiency.

4.4 LONG VIDEO GENERATION EVALUATION

We evaluated our long video generation framework using both simple and complex prompts. As shown in Figure 8, our model maintains consistent perceptual quality over time, whereas baseline methods suffer from rapid quality decay due to error accumulation. A large-scale user study (93 participants, 150 comparisons per model pair) on the first 50 prompts from MovieGenBench further validated

| Attention Window Size | Total Score | In-Flight Latency (s) | Memory (GB) |
|-----------------------|-------------|-----------------------|-------------|
| 3 | 80.9 | 0.37 ± 0.01 | 14.9 |
| 6 | 81.3 | 0.38 ± 0.01 | 15.8 |
| 9 | 82.3 | 0.40 ± 0.01 | 16.6 |
| 12 | 84.3 | 0.46 ± 0.02 | 17.5 |
| 15 | 84.2 | 0.51 ± 0.02 | 18.4 |
| 18 | 84.4 | 0.54 ± 0.02 | 19.2 |
| 21 | 84.5 | 0.59 ± 0.02 | 20.1 |
| 24 | 84.3 | 0.64 ± 0.02 | 20.9 |
| 27 | 84.5 | 0.68 ± 0.02 | 21.8 |

Table 4: KV cache scaling analysis



Figure 9: Illustration of long video generation with dynamic prompting. This feature allows for the integration of new prompts at arbitrary time points, facilitating the creation of coherent long videos with changing narratives. The specific prompts used for each segment are detailed in the appendix.

our method’s superiority in overall visual quality, text faithfulness, and long-term consistency. User study results, consistent with the qualitative comparison in Figure 7, confirm that baseline methods degrade with issues like saturation distortion, while our approach sustains high quality. A key feature of our framework is its support for *dynamic prompting* (Figure 9), allowing users to input new text descriptions at any timeline point to create complex narratives with evolving scenes and actions.

5 CONCLUDING REMARKS

In this work, we introduce Diagonal Distillation, a novel framework for efficient autoregressive video generation. It leverages temporal dependencies across video chunks and denoising steps through an asymmetric denoising strategy—allocating more steps to early chunks and progressively fewer to later ones. This design significantly reduces the total number of denoising steps while preserving motion coherence and visual quality. Diagonal Forcing explicitly models the temporal denoising trajectory, reducing error accumulation and aligning training with inference for stable long-range synthesis. Additionally, Flow Distribution Matching ensures dynamic consistency under strict step constraints by aligning the optical flow distributions of generated and real videos. Extensive experiments demonstrate our method’s superior trade-off between efficiency and quality.

540 **ETHICS STATEMENT**

541

542 While the real-time video generation technology presented in this study significantly improves
 543 generation efficiency (achieving a 277.3x speedup compared to the baseline model), we are fully
 544 aware of its dual-use nature. This technology could potentially be misused to create misleading
 545 content or deepfake videos. To mitigate this risk, we commit to embedding usage guidelines and
 546 restrictions when open-sourcing the code and models, and we advocate for the adoption of traceability
 547 technologies such as digital watermarks and content authentication. Concurrently, this technology
 548 holds significant positive potential in fields such as education, the creative industries, and assistive
 549 tools. We aim to maximize its societal benefits and minimize potential harms through ongoing
 550 discussions on technology ethics and responsible release practices.

551

552 **REPRODUCIBILITY STATEMENT**

553

554 For detailed reproducibility information, including full implementation details, training configurations,
 555 hyperparameters, and evaluation protocols, please refer to the appendix sections. All source code,
 556 trained model weights, and configuration files will be released to ensure the full reproducibility of
 557 our results.

559 **REFERENCES**

560

- 561 Marianne Arriola, Aaron Gokaslan, Justin T Chiu, Zhihan Yang, Zhixuan Qi, Jiaqi Han, Sub-
 562 ham Sekhar Sahoo, and Volodymyr Kuleshov. Block diffusion: Interpolating between autoregres-
 563 sive and diffusion language models. *arXiv preprint arXiv:2503.09573*, 2025. 2, 3
- 564 Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik
 565 Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion: Scaling
 566 latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023a. 1, 3
- 567 Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler,
 568 and Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion
 569 models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*,
 570 pp. 22563–22575, 2023b. 1, 3
- 571 Ollin Boer Bohan. Taehv: Tiny autoencoder for hunyuan video. <https://github.com/madebyollin/taehv>, 2025. 7
- 572 Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe
 573 Taylor, Troy Luhman, Eric Luhman, et al. Video generation models as world simulators. *OpenAI*
 574 *Blog*, 1(8):1, 2024. 1, 3
- 575 Jake Bruce, Michael D Dennis, Ashley Edwards, Jack Parker-Holder, Yuge Shi, Edward Hughes,
 576 Matthew Lai, Aditi Mavalankar, Richie Steigerwald, Chris Apps, et al. Genie: Generative
 577 interactive environments. In *Forty-first International Conference on Machine Learning*, 2024. 1, 3
- 578 Boyuan Chen, Diego Martí Monsó, Yilun Du, Max Simchowitz, Russ Tedrake, and Vincent Sitzmann.
 579 Diffusion forcing: Next-token prediction meets full-sequence diffusion. *Advances in Neural*
 580 *Information Processing Systems*, 37:24081–24125, 2024. 2, 3
- 581 Guibin Chen, Dixuan Lin, Jiangping Yang, Chunze Lin, Junchen Zhu, Mingyuan Fan, Hao Zhang,
 582 Sheng Chen, Zheng Chen, Chengcheng Ma, et al. Skyreels-v2: Infinite-length film generative
 583 model. *arXiv preprint arXiv:2504.13074*, 2025. 7
- 584 Haoge Deng, Ting Pan, Haiwen Diao, Zhengxiong Luo, Yufeng Cui, Huchuan Lu, Shiguang Shan,
 585 Yonggang Qi, and Xinlong Wang. Autoregressive video generation without vector quantization.
 586 *arXiv preprint arXiv:2412.14169*, 2024. 1, 3
- 587 Kaifeng Gao, Jiaxin Shi, Hanwang Zhang, Chunping Wang, and Jun Xiao. Vid-gpt: Introducing
 588 gpt-style autoregressive generation in video diffusion models. *arXiv preprint arXiv:2406.10981*,
 589 2024a. 1

- 594 Kaifeng Gao, Jiaxin Shi, Hanwang Zhang, Chunping Wang, Jun Xiao, and Long Chen. Ca2-vdm:
 595 Efficient autoregressive video diffusion model with causal generation and cache sharing. *arXiv*
 596 preprint arXiv:2411.16375, 2024b. 2, 3
- 597 Songwei Ge, Thomas Hayes, Harry Yang, Xi Yin, Guan Pang, David Jacobs, Jia-Bin Huang, and
 598 Devi Parikh. Long video generation with time-agnostic vqgan and time-sensitive transformer. In
 599 *European Conference on Computer Vision*, pp. 102–118. Springer, 2022. 1
- 600 Yuchao Gu, Weijia Mao, and Mike Zheng Shou. Long-context autoregressive video modeling with
 601 next-frame prediction. *arXiv preprint arXiv:2503.19325*, 2025. 2, 3
- 602 Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P
 603 Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High definition
 604 video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022. 3
- 605 Jinyi Hu, Shengding Hu, Yuxuan Song, Yufei Huang, Mingxuan Wang, Hao Zhou, Zhiyuan Liu,
 606 Wei-Ying Ma, and Maosong Sun. Acdit: Interpolating autoregressive conditional modeling and
 607 diffusion transformer. *arXiv preprint arXiv:2412.07720*, 2024. 2, 3
- 608 Xun Huang, Zhengqi Li, Guande He, Mingyuan Zhou, and Eli Shechtman. Self forcing: Bridging the
 609 train-test gap in autoregressive video diffusion. *arXiv preprint arXiv:2506.08009*, 2025. 1, 2, 3, 7
- 610 Ziqi Huang, Yinan He, Jiashuo Yu, Fan Zhang, Chenyang Si, Yuming Jiang, Yuanhan Zhang, Tianxing
 611 Wu, Qingyang Jin, Nattapol Chanpaisit, et al. Vbench: Comprehensive benchmark suite for video
 612 generative models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern
 613 Recognition*, pp. 21807–21818, 2024. 7
- 614 Yang Jin, Zhicheng Sun, Ningyuan Li, Kun Xu, Hao Jiang, Nan Zhuang, Quzhe Huang, Yang Song,
 615 Yadong Mu, and Zhouchen Lin. Pyramidal flow matching for efficient video generative modeling.
 616 *arXiv preprint arXiv:2410.05954*, 2024. 1, 3
- 617 A Jolicoeur-Martineau. The relativistic discriminator: A key element missing from standard gan.
 618 *arXiv preprint arXiv:1807.00734*, 2018. 1
- 619 Diederik Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models. *Advances
 620 in neural information processing systems*, 34:21696–21707, 2021. 3
- 621 Dan Kondratyuk, Lijun Yu, Xiuye Gu, José Lezama, Jonathan Huang, Grant Schindler, Rachel
 622 Hornung, Vighnesh Birodkar, Jimmy Yan, Ming-Chang Chiu, et al. Videopoet: A large language
 623 model for zero-shot video generation. *arXiv preprint arXiv:2312.14125*, 2023. 1, 3
- 624 Weijie Kong, Qi Tian, Zijian Zhang, Rox Min, Zuozhuo Dai, Jin Zhou, Jiangfeng Xiong, Xin Li,
 625 Bo Wu, Jianwei Zhang, et al. Hunyanvideo: A systematic framework for large video generative
 626 models. *arXiv preprint arXiv:2412.03603*, 2024. 1, 3
- 627 Alex M Lamb, Anirudh Goyal ALIAS PARTH GOYAL, Ying Zhang, Saizheng Zhang, Aaron C
 628 Courville, and Yoshua Bengio. Professor forcing: A new algorithm for training recurrent networks.
 629 *Advances in neural information processing systems*, 29, 2016. 3
- 630 Qing Li, Xun Tang, Junkun Peng, Yuanzheng Tan, and Yong Jiang. Latency reducing in real-time
 631 internet video transport: A survey. Available at SSRN 4654242, 2023. 3
- 632 Tianhong Li, Yonglong Tian, He Li, Mingyang Deng, and Kaiming He. Autoregressive image
 633 generation without vector quantization. *Advances in Neural Information Processing Systems*, 37:
 634 56424–56445, 2024a. 3
- 635 Zhengqi Li, Qianqian Wang, Noah Snavely, and Angjoo Kanazawa. Infinitenature-zero: Learning
 636 perpetual view generation of natural scenes from single images. In *European conference on
 637 computer vision*, pp. 515–534. Springer, 2022. 3
- 638 Zongyi Li, Shujie Hu, Shujie Liu, Long Zhou, Jeongsoo Choi, Lingwei Meng, Xun Guo, Jinyu Li,
 639 Hefei Ling, and Furu Wei. Arlon: Boosting diffusion transformers with autoregressive models for
 640 long video generation. *arXiv preprint arXiv:2410.20502*, 2024b. 2, 3

- 648 Shanchuan Lin, Anran Wang, and Xiao Yang. Sdxl-lightning: Progressive adversarial diffusion
 649 distillation. *arXiv preprint arXiv:2402.13929*, 2024. 3
 650
- 651 Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching
 652 for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022. 7
 653
- 654 Haozhe Liu, Shikun Liu, Zijian Zhou, Mengmeng Xu, Yanping Xie, Xiao Han, Juan C Pérez, Ding
 655 Liu, Kumara Kahatapitiya, Menglin Jia, et al. Mardini: Masked autoregressive diffusion for video
 656 generation at scale. *arXiv preprint arXiv:2410.20280*, 2024a. 3
 657
- 658 Jinxiu Liu, Shaoheng Lin, Yinxiao Li, and Ming-Hsuan Yang. Dynamicscaler: Seamless and scalable
 659 video generation for panoramic scenes. In *Proceedings of the Computer Vision and Pattern
 Recognition Conference*, pp. 6144–6153, 2025. 3
 660
- 661 Yaofang Liu, Yumeng Ren, Xiaodong Cun, Aitor Artola, Yang Liu, Tieyong Zeng, Raymond H Chan,
 662 and Jean-michel Morel. Redefining temporal modeling in video diffusion: The vectorized timestep
 663 approach. *arXiv preprint arXiv:2410.03160*, 2024b. 2, 3
 664
- 665 Zhijun Liu, Shuai Wang, Sho Inoue, Qibing Bai, and Haizhou Li. Autoregressive diffusion transformer
 666 for text-to-speech synthesis. *arXiv preprint arXiv:2406.05551*, 2024c. 2, 3
 667
- 668 Weijian Luo, Zemin Huang, Zhengyang Geng, J Zico Kolter, and Guo-jun Qi. One-step diffusion
 669 distillation through score implicit matching. *Advances in Neural Information Processing Systems*,
 37:115377–115408, 2024. 3
 670
- 671 William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of
 the IEEE/CVF international conference on computer vision*, pp. 4195–4205, 2023. 1
 672
- 673 Adam Polyak, Amit Zohar, Andrew Brown, Andros Tjandra, Animesh Sinha, Ann Lee, Apoorv Vyas,
 674 Bowen Shi, Chih-Yao Ma, Ching-Yao Chuang, et al. Movie gen: A cast of media foundation
 675 models. *arXiv preprint arXiv:2410.13720*, 2024. 1, 3
 676
- 677 Shuhuai Ren, Shuming Ma, Xu Sun, and Furu Wei. Next block prediction: Video generation via
 678 semi-autoregressive modeling. *arXiv preprint arXiv:2502.07737*, 2025. 1, 3
 679
- 680 Yuxi Ren, Xin Xia, Yanzuo Lu, Jiacheng Zhang, Jie Wu, Pan Xie, Xing Wang, and Xuefeng Xiao.
 681 Hyper-sd: Trajectory segmented consistency model for efficient image synthesis. *Advances in
 682 Neural Information Processing Systems*, 37:117340–117362, 2024. 3
 683
- 684 Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. *arXiv
 preprint arXiv:2202.00512*, 2022. 3
 685
- 686 Axel Sauer, Frederic Boesel, Tim Dockhorn, Andreas Blattmann, Patrick Esser, and Robin Rombach.
 687 Fast high-resolution image synthesis with latent adversarial diffusion distillation. In *SIGGRAPH
 Asia 2024 Conference Papers*, pp. 1–11, 2024a. 3
 688
- 689 Axel Sauer, Dominik Lorenz, Andreas Blattmann, and Robin Rombach. Adversarial diffusion
 690 distillation. In *European Conference on Computer Vision*, pp. 87–103. Springer, 2024b. 3
 691
- 692 Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models, 2023. URL
<https://arxiv.org/abs/2303.01469>. 1, 3
 693
- 694 Hansi Teng, Hongyu Jia, Lei Sun, Lingzhi Li, Maolin Li, Mingqiu Tang, Shuai Han, Tianning Zhang,
 695 WQ Zhang, Weifeng Luo, et al. Magi-1: Autoregressive video generation at scale. *arXiv preprint
 arXiv:2505.13211*, 2025. 1, 7, 9
 696
- 697 Ruben Villegas, Mohammad Babaeizadeh, Pieter-Jan Kindermans, Hernan Moraldo, Han Zhang,
 698 Mohammad Taghi Saffar, Santiago Castro, Julius Kunze, and Dumitru Erhan. Phenaki: Variable
 699 length video generation from open domain textual description. *arXiv preprint arXiv:2210.02399*,
 700 2022. 1, 3
 701
- Carl Vondrick, Hamed Pirsiavash, and Antonio Torralba. Generating videos with scene dynamics.
Advances in neural information processing systems, 29, 2016. 1

- 702 Team Wan, Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, Chen-Wei Xie, Di Chen, Feiwu Yu,
 703 Haiming Zhao, Jianxiao Yang, et al. Wan: Open and advanced large-scale video generative models.
 704 *arXiv preprint arXiv:2503.20314*, 2025. 1, 3, 7
 705
- 706 Wenhao Wang and Yi Yang. Vidprom: A million-scale real prompt-gallery dataset for text-to-video
 707 diffusion models. *Advances in Neural Information Processing Systems*, 37:65618–65642, 2024. 7
 708
- 709 Yuqing Wang, Tianwei Xiong, Daquan Zhou, Zhijie Lin, Yang Zhao, Bingyi Kang, Jiashi Feng, and
 710 Xihui Liu. Loong: Generating minute-level long videos with autoregressive language models.
 711 *arXiv preprint arXiv:2410.02757*, 2024. 1, 3
 712
- 713 Zhengyi Wang, Cheng Lu, Yikai Wang, Fan Bao, Chongxuan Li, Hang Su, and Jun Zhu. Pro-
 714 lifedreamer: High-fidelity and diverse text-to-3d generation with variational score distillation.
 715 *Advances in neural information processing systems*, 36:8406–8441, 2023. 1
 716
- 717 Dirk Weissenborn, Oscar Täckström, and Jakob Uszkoreit. Scaling autoregressive video models.
 718 *arXiv preprint arXiv:1906.02634*, 2019. 1, 3
 719
- 720 Wenming Weng, Ruoyu Feng, Yanhui Wang, Qi Dai, Chunyu Wang, Dacheng Yin, Zhiyuan Zhao,
 721 Kai Qiu, Jianmin Bao, Yuhui Yuan, et al. Art-v: Auto-regressive text-to-video generation with
 722 diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern
 Recognition*, pp. 7395–7405, 2024. 1, 2, 3
 723
- 724 Wilson Yan, Yunzhi Zhang, Pieter Abbeel, and Aravind Srinivas. Videogpt: Video generation using
 725 vq-vae and transformers. *arXiv preprint arXiv:2104.10157*, 2021. 1, 3
 726
- 727 Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming Yang,
 728 Wenyi Hong, Xiaohan Zhang, Guanyu Feng, et al. Cogvideox: Text-to-video diffusion models
 729 with an expert transformer. *arXiv preprint arXiv:2408.06072*, 2024. 1, 3
 730
- 731 Tianwei Yin, Michaël Gharbi, Taesung Park, Richard Zhang, Eli Shechtman, Fredo Durand, and Bill
 732 Freeman. Improved distribution matching distillation for fast image synthesis. *Advances in neural
 733 information processing systems*, 37:47455–47487, 2024a. 3
 734
- 735 Tianwei Yin, Michaël Gharbi, Richard Zhang, Eli Shechtman, Fredo Durand, William T Freeman,
 736 and Taesung Park. One-step diffusion with distribution matching distillation. In *Proceedings of the
 737 IEEE/CVF conference on computer vision and pattern recognition*, pp. 6613–6623, 2024b. 1, 3
 738
- 739 Tianwei Yin, Qiang Zhang, Richard Zhang, William T Freeman, Fredo Durand, Eli Shechtman, and
 740 Xun Huang. From slow bidirectional to fast autoregressive video diffusion models. In *Proceedings
 741 of the Computer Vision and Pattern Recognition Conference*, pp. 22963–22974, 2025. 1, 2, 3, 7
 742
- 743 Tianyuan Zhang, Sai Bi, Yicong Hong, Kai Zhang, Fujun Luan, Songlin Yang, Kalyan Sunkavalli,
 744 William T Freeman, and Hao Tan. Test-time training done right. *arXiv preprint arXiv:2505.23884*,
 745 2025a. 2, 3
 746
- 747 Yuan Zhang, Jiacheng Jiang, Guoqing Ma, Zhiying Lu, Haoyang Huang, Jianlong Yuan, and
 748 Nan Duan. Generative pre-trained autoregressive diffusion transformer. *arXiv preprint
 749 arXiv:2505.07344*, 2025b. 2, 3
 750
- 751
- 752
- 753
- 754
- 755