

A Spatial Feature Extraction

We initially utilize a CRAFT OCR [4] to outline the text area (*TA*) of the frame F_d . Then we adopt the prompt encoder heuristic in the Segment Anything Model (SAM) [18] to convert *TA* into prompt embeddings X_{vt}^d to obtain the feature vector of the location and content information of the text area. At the same time, we input F_d into the pre-trained ViT to generate the initial encoding X_v^d to capture the spatial features of the image in a single frame.

For the obtained X_{vt}^d and X_v^d , we optimize the representation of visual features X_v^d through a bidirectional attention mechanism, making it more focused on the relative representation of subtitle features in space. The process mainly takes the following steps:

We input X_{vt}^d and X_v^d into a bidirectional attention block. The bidirectional attention block first performs a self-attention (Self-Att) operation on X_{vt}^d to enhance its own feature representation.

$$\hat{X}_{vt}^d = \frac{(X_{vt}^d + \text{Self-Att}(X_{vt}^d)) - \mu}{\sqrt{\sigma^2 + \epsilon}} \sim \text{NORM}(X_{vt}^d + \text{Self-Att}(X_{vt}^d)) \quad (16)$$

where Self-Att is self-attention, and we normalize all features within the sample to have zero mean and unit variance to help alleviate the problem of internal covariate shift. Later, we use **NORM** to simplify this normalization process. Then, we perform cross-attention (Cross-Att) operations on the optimized features \hat{X}_{vt}^d and X_v^d to interact with each other and further optimize the \hat{X}_{vt}^d .

$$\dot{X}_{vt}^d = \text{NORM}(\hat{X}_{vt}^d + \text{Cross-Att}(\hat{X}_{vt}^d, X_v^d)) \quad (17)$$

Lastly, a cross-attention operation is performed on X_v^d to make it interact with the updated \dot{X}_{vt}^d from MLP to obtain a more accurate visual feature representation.

$$\hat{X}_v^d = \text{NORM}(X_v^d + \text{Cross-Att}(\text{MLP}(\dot{X}_{vt}^d), X_v^d)) \quad (18)$$

After the attention process, the updated \hat{X}_v^d is down-sampled through two convolutional layers and then expanded to obtain the spatial pattern features of a single frame X^d . Finally, we use self-attention to fuse these single-frame spatial features and obtain a unified spatial feature representation through average pooling (MEAN).

$$X_{spa} = \text{MEAN}(\text{Self-Att}([X^1, X^2, \dots, X^d])) \quad (19)$$

B Domain datasets

Table 5: Domains division of FakeSV and FakeTT datasets.

Domain	FakeSV	FakeTT
Culture	278	325
Disaster	1836	182
Education	195	28
Finance	116	9
Health	841	121
Politics	296	715
Military	121	84
Science	234	271
Society	1621	337

We use the Qwen2.5-72B [41] combined with manual verification for label classification based on the text and fakesv's video collection introduction [28]. We divide the above two datasets into nine domains, respectively. Specific domains' names and numbers are detailed and shown in Table 5.

The specific prompt template is: “*Here are some fields: Society, Health, Disaster, Culture, Education, Finance, Politics, Science, Military. And some examples from different domains [Some examples (Society: [case text], Health: [case text], ...)]. You can then compare and analyze the [input text] with the examples above to choose the most suitable Domain.*”

We divide the above two datasets into nine domains, respectively. Specific domains' names and numbers are detailed and shown in Table 5. Since the number of some domains is small (Education, Finance, and Military in the FakeTT dataset), we conduct merging experiments with them.

C More Experimental Analyses

C.1 Training Strategy

We define the total loss function $Loss_{final}$ as the combination of the **Cross-modal Interpolation Distillation** loss $Loss_{dis}$ and the **Cross-modal Invariance Fusion** loss $Loss_{es}$ (emotion and sentiment features) and $Loss_{st}$ (spatial and temporal features). The specifically $Loss_{final}$ we define is as follows:

$$Loss_{final} = \alpha Loss_{es} + \beta Loss_{st} + \gamma Loss_{dis}, \quad (20)$$

where α , β , and γ denote the weighting factors.

C.2 Experimental Implementations

Our experiments include two aspects: domain generalization experiments and general domain experiments. The comparison methods selected in the domain generalization experiments are mainly the best performing and latest methods selected from general domain experiments for comparison. All our experiments were conducted on an NVIDIA RTX A6000 (48G). The α in the loss functions $Loss_{final}$ is 0.1, β is 3, and γ is 0.05. We choose the learning rates of 5e-5 for FakeSV and 1e-3 for FakeTT.

For more details on the DOCTOR model, our anonymous source code is available at <https://anonymous.4open.science/r/DOCTOR-D1D2>.

C.3 Comparison of Domain Generalization Star Map

As shown in Figure 6, the experiment present a comparative analysis of six misinformation detection methods (FakingRecipe, SV-FEND, OpEvFake, CMRF*, MMDFND*, and DOCTOR) across two distinct datasets: FakeSV and FakeTT. The radar visualization effectively illustrates each method's performance across multiple content categories, including Disaster, Society, Health, Culture, Politics, Science, Education, Finance, and Military for FakeSV, with a combined Education/Finance/Military category for FakeTT. The charts reveal DOCTOR consistently outperforming other methods across most categories in both datasets, particularly excelling in Military (93.65%) for FakeSV and Culture (88.31%) for FakeTT.

1277

1278

1279

1280

1281

1282

1283

1284

1285

1286

1287

1288

1289

1290

1291

1292

1293

1294

1295

1296

1297

1298

1299

1300

1301

1302

1303

1304

1305

1306

1307

1308

1309

1310

1311

1312

1313

1314

1315

1316

1317

1318

1319

1320

1321

1322

1323

1324

1325

1326

1327

1328

1329

1330

1331

1332

1333

1334

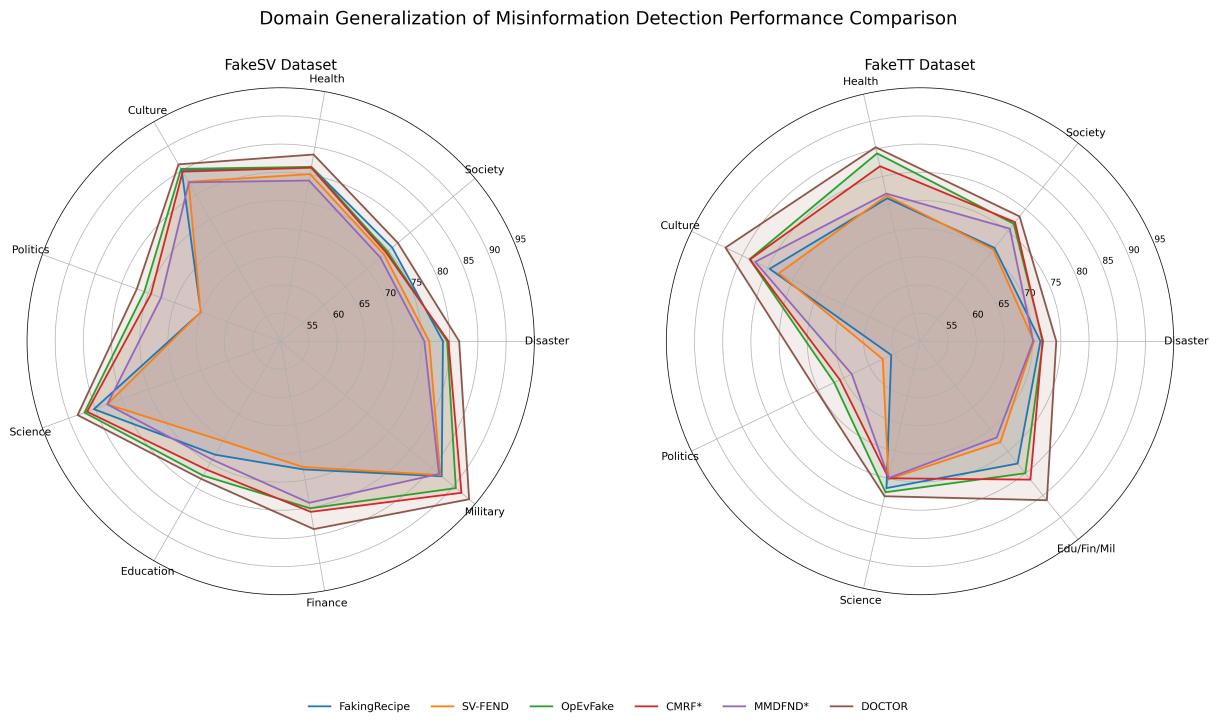


Figure 6: Comparative analysis of six misinformation detection methods (FakingRecipe, SV-FEND, OpEvFake, CMRF*, MMDFND*, and DOCTOR) across two distinct datasets: FakeSV and FakeTT. * represents the domain generalization method.

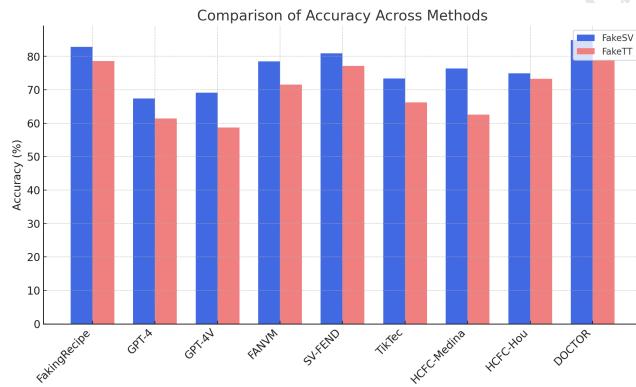


Figure 7: Accuracy comparison of different methods on FakeSV and FakeTT datasets.

C.4 More comparative accuracy analysis

The above two figures of Figure 7 and Figure 8 show the comparison of the Accuracy and F1-score of different methods on FakeSV and FakeTT datasets. It can be seen that the DOCTOR method achieved the best performance on both datasets.

The above three figures of Figure 11, Figure 10, and Figure 9 show the comparison of the Precision, Recall, and F1-score of different methods on the Fake category on the FakeSV and FakeTT datasets. It can be seen that the DOCTOR method performs outstandingly in Recall and F1-score on the Fake category.

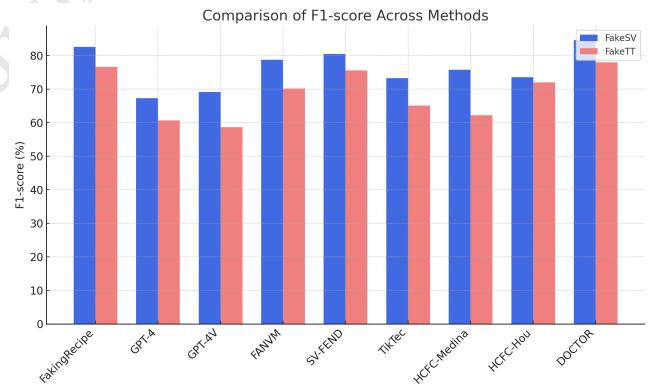


Figure 8: F1 comparison of different methods on FakeSV and FakeTT datasets.

The above three figures of Figure 14, Figure 13, and Figure 12 show the comparison of the Precision, Recall, and F1-score of the Real category of different methods on the FakeSV and FakeTT datasets. From the results, the DOCTOR method performs best in Precision while maintaining a high level in Recall and F1-score, and the overall performance is stable.

C.5 Discussion on DOCTOR Model Efficiency

Despite the advanced features and complex mechanisms integrated into the DOCTOR model, as shown in Figure 15, our experiments

1335

1336

1337

1338

1339

1340

1341

1342

1343

1344

1345

1346

1347

1348

1349

1350

1351

1352

1353

1354

1355

1356

1357

1358

1359

1360

1361

1362

1363

1364

1365

1366

1367

1368

1369

1370

1371

1372

1373

1374

1375

1376

1377

1378

1379

1380

1381

1382

1383

1384

1385

1386

1387

1388

1389

1390

1391

1392

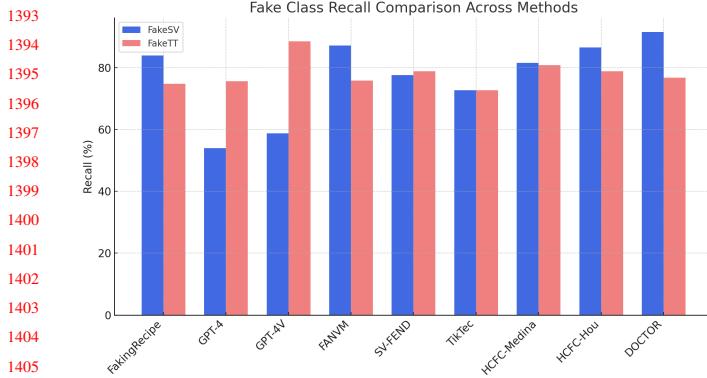


Figure 9: Recall comparison of different methods on FakeSV and FakeTT datasets (Fake class).

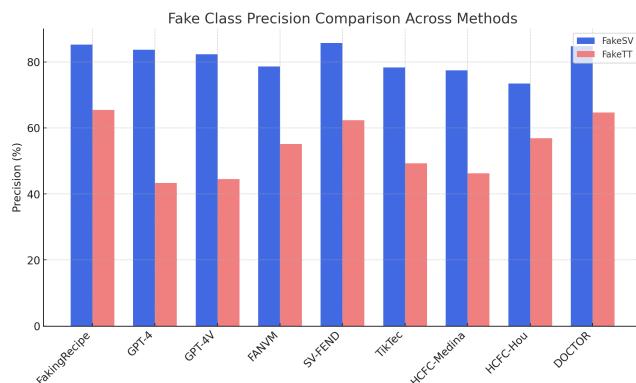


Figure 10: Precision comparison of different methods on FakeSV and FakeTT datasets (Fake class).

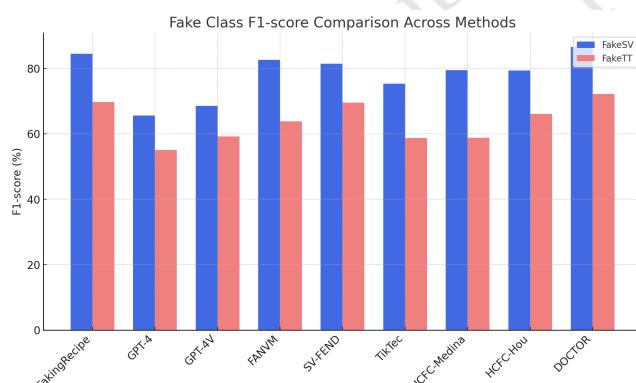


Figure 11: F1 comparison of different methods on FakeSV and FakeTT datasets (Fake class).

reveal that the computational overhead does not experience a significant increase. The incorporation of cross-modal interpolation distillation and multi-modal invariance fusion, while enhancing the model's generalization ability across different domains, does not lead

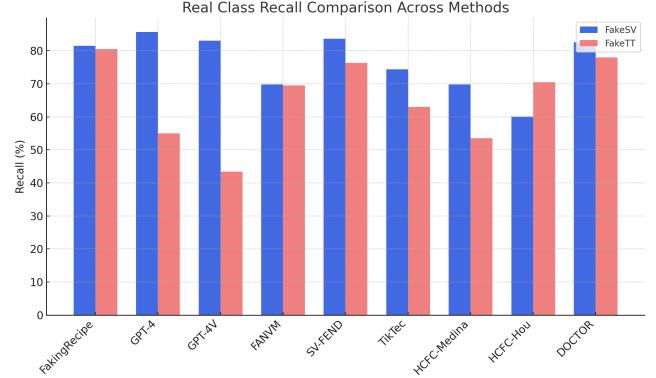


Figure 12: Recall comparison of different methods on FakeSV and FakeTT datasets (Real class).

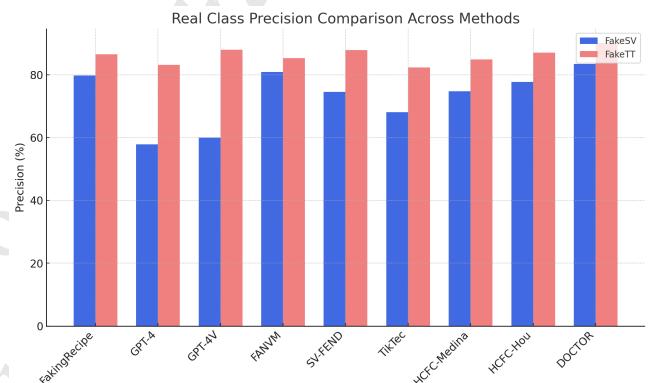


Figure 13: Precision comparison of different methods on FakeSV and FakeTT datasets (Real class).

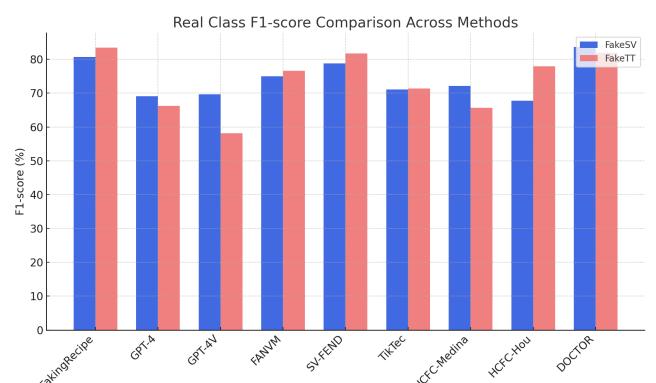


Figure 14: F1 comparison of different methods on FakeSV and FakeTT datasets (Real class).

to a substantial rise in computational cost. This indicates that DOCTOR is able to maintain efficient performance, even with the addition of multiple modalities and the complexity of noise-based feature refinement. The efficiency of DOCTOR is crucial for real-world

1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403
1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457
1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508

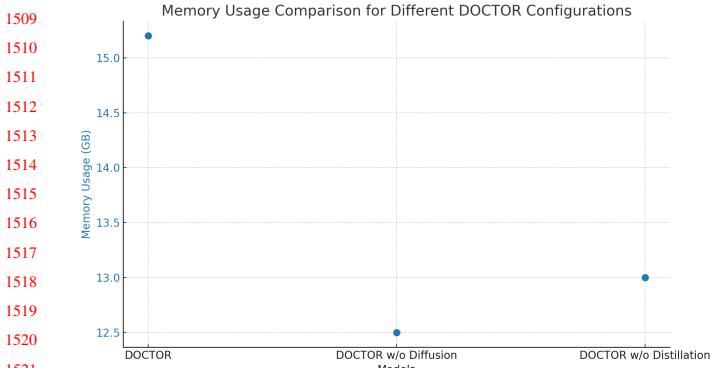


Figure 15: Discussion on DOCTOR Model Efficiency.

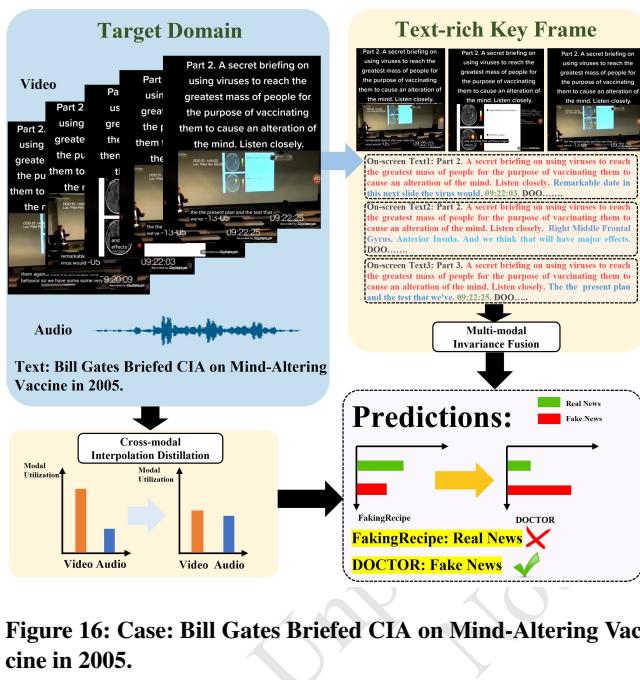


Figure 16: Case: Bill Gates Briefed CIA on Mind-Altering Vaccine in 2005.

applications, where both accuracy and computational efficiency are paramount.

C.6 Case Study

As shown in Figure 16. Case: Bill Gates Briefed CIA on Mind-Altering Vaccine in 2005.

In this case, the DOCTOR model demonstrates strong domain generalization capabilities by correctly identifying the video as misinformation, while traditional methods like FakingRecipe misclassify it as truth. They struggle to adapt when encountering unseen distributions, variations in video styles, or novel linguistic manipulations. In contrast, DOCTOR incorporates multi-modal invariance fusion and cross-modal interpolation distillation, allowing it to generalize across different domains by capturing latent relationships between text, audio, and video. This enables it to detect cross-modal

inconsistencies, making it more robust against out-of-domain misinformation scenarios.

D Temporal Position Encoding

Temporal Position Encoding (TPE) is designed to capture the coarse temporal location of each segment within the entire video. To achieve this, we divide the video into B equal-frequency temporal bins and assign each segment to one of them based on its position in the sequence. Each bin is associated with a learnable embedding vector, which is used to represent the temporal context of segments assigned to that bin. This encoding helps the model understand whether a segment is located at the beginning, middle, or end of the video, which can be crucial for modeling temporal editing patterns.

Algorithm 2 Temporal Position Encoding (TPE)

Require: Number of segments n , number of bins B , embedding dimension d
Ensure: Temporal position embeddings $\{TPE^1, \dots, TPE^n\}$

- 1: $s \leftarrow \lfloor \frac{n}{B} \rfloor$
- 2: Initialize embedding matrix $E \in \mathbb{R}^{B \times d}$
- 3: **for** $i = 1$ to n **do**
- 4: $b_i \leftarrow \min \left(\lfloor \frac{i-1}{s} \rfloor, B-1 \right)$
- 5: $TPE^i \leftarrow E[b_i]$
- 6: **end for**
- 7: **return** $\{TPE^1, \dots, TPE^n\}$

E Duration Encoding

Duration Encoding (DE) encodes the fine-grained temporal position of each segment using a fixed sinusoidal scheme, similar to the original positional encoding used in Transformers. This method ensures that each segment has a unique continuous representation based on its position index. For each segment i , and for each embedding dimension j , the encoding alternates between sine and cosine functions of different frequencies. This encoding provides precise temporal cues, enabling the model to capture temporal dependencies and alignment across modalities.

Algorithm 3 Duration Encoding (DE)

Require: Number of segments n , encoding dimension d
Ensure: Duration encodings $\{DE^1, \dots, DE^n\}$

- 1: Initialize $DE \in \mathbb{R}^{n \times d}$ with zeros
- 2: **for** $i = 0$ to $n - 1$ **do**
- 3: **for** $j = 0$ to $d - 1$ **do**
- 4: $\theta \leftarrow \frac{i}{2 \cdot \lfloor j/2 \rfloor} \cdot \frac{10000}{d}$
- 5: **if** j is even **then**
- 6: $DE^{i+1}[j] \leftarrow \sin(\theta)$
- 7: **else**
- 8: $DE^{i+1}[j] \leftarrow \cos(\theta)$
- 9: **end if**
- 10: **end for**
- 11: **end for**
- 12: **return** $\{DE^1, \dots, DE^n\}$

1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619
1620
1621
1622
1623

1625 F Predication

1626

1627 Algorithm 4 Prediction via Feature Diffusion and Distillation Fusion

1628

```

1: Input: Time feature  $X_{\text{time}}$ , Spatial feature  $X_{\text{spa}}$ ,
   Semantic feature  $X_{\text{sem}}$ , Emotional feature  $X_{\text{emo}}$ ,
   Video feature  $X_{\text{video}}$ , Audio feature  $X_{\text{audio}}$ 
2: Output: Final prediction score  $\hat{Y}_{\text{final}}$ 
3: {Step 1: Feature fusion and diffusion}
4:  $X_{\text{all}} \leftarrow \text{Concatenate}(X_{\text{time}}, X_{\text{spa}}, X_{\text{sem}}, X_{\text{emo}})$ 
5:  $X_{\text{diffused}} \leftarrow \text{ApplyDiffusion}(X_{\text{all}})$ 
6: {Step 2: Time and space prediction branch}
7:  $X_{\text{ts}} \leftarrow \text{Extract}(X_{\text{diffused}}, [\text{time}, \text{spa}])$ 
8:  $\hat{Y}_{\text{st}} \leftarrow \text{MLP}_{\text{TS}}(X_{\text{ts}})$ 
9: {Step 3: Semantic and emotional prediction branch}
10:  $X_{\text{se}} \leftarrow \text{Extract}(X_{\text{diffused}}, [\text{sem}, \text{emo}])$ 
11:  $\hat{Y}_{\text{es}} \leftarrow \text{MLP}_{\text{SE}}(X_{\text{se}})$ 
12: {Step 4: Video-Audio distillation output}
13:  $\hat{Y}_{\text{dis}} \leftarrow \text{Distill}(X_{\text{video}}, X_{\text{audio}})$ 
14: {Step 5: Late fusion to get final prediction}
15:  $\hat{Y}_{\text{fusion}} \leftarrow \hat{Y}_{\text{st}} \times \tanh(\hat{Y}_{\text{es}})$ 
16:  $\hat{Y}_{\text{final}} \leftarrow \hat{Y}_{\text{fusion}} \times \hat{Y}_{\text{dis}}$ 
17: return  $\hat{Y}_{\text{final}}$ 
```

1647

1648

1649

1650

1651

1652

1653

1654

1655

1656

1657

1658

1659

1660

1661

1662

1663

1664

1665

1666

1667

1668

1669

1670

1671

1672

1673

1674

1675

1676

1677

1678

1679

1680

1681

1682

To generate the final prediction, we design a multi-branch architecture that integrates temporal, spatial, semantic, and emotional features through a diffusion-enhanced fusion mechanism. As illustrated in Algorithm 4, the input features X_{time} , X_{spa} , X_{sem} , and X_{emo} are first concatenated and passed through a diffusion process to enhance inter-feature dependencies and latent structure alignment. The diffused representation X_{diffused} is then split into two branches:

Time-Spatial Branch: We extract the time and spatial-related features to form X_{ts} , which is passed through a dedicated MLP_{TS} to yield an intermediate prediction \hat{Y}_{st} .

Semantic-Emotional Branch: Similarly, semantic and emotional components are extracted as X_{se} , which is processed via another MLP_{SE} to produce \hat{Y}_{es} .

In parallel, we distill information from both video and audio modalities using a specialized distillation module, resulting in \hat{Y}_{dis} .

For final prediction, we adopt a late fusion strategy where the time-spatial output is modulated by the nonlinear activation of the semantic-emotional output: $\hat{Y}_{\text{fusion}} = \hat{Y}_{\text{st}} \cdot \tanh(\hat{Y}_{\text{es}})$. This fusion is then further modulated by the distilled multi-modal prediction: $\hat{Y}_{\text{final}} = \hat{Y}_{\text{fusion}} \cdot \hat{Y}_{\text{dis}}$. This design effectively captures complementary information from all feature domains and ensures robustness through multi-modal consistency.

For specific code implementation, please refer to the code link.

1683

1684

1685

1686

1687

1688

1689

1690

1691

1692

1693

1694

1695

1696

1697

1698

1699

1700

1701

1702

1703

1704

1705

1706

1707

1708

1709

1710

1711

1712

1713

1714

1715

1716

1717

1718

1719

1720

1721

1722

1723

1724

1725

1726

1727

1728

1729

1730

1731

1732

1733

1734

1735

1736

1737

1738

1739