# Appendix of "AutoAD: Movie Description in Context"

We first show the details of the AD collection pipeline (Sec. A) with qualitative text examples (Sec. B). Then we describe additional implementation details (Sec. D) with extra qualitative movie AD examples (Sec. E). Finally, we list the movie IDs used in our MAD-v2 split (Sec. F).
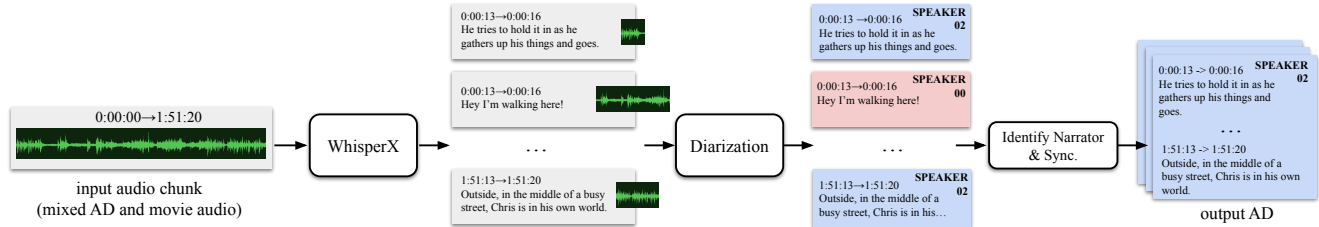
## A. Additional Details of the AD Collection Pipeline



Figure A.1. **A schematic of our AD collection pipeline.** The pipeline takes the audio file (with mixed AD and movie audio) as input, and automatically outputs the AD in the text form with the corresponding timestamps.

### A.1. AD Collection Pipeline for MAD-v2 dataset

Collecting movie AD has two main challenges. First, in the audio files (*e.g.* from AudioVault) the movie AD is *fused* with the original movie audio, *i.e.* on the same audio track. The pipeline needs to identify the AD speaker among the movie characters accurately. Second, for the same movie, the audio files from AudioVault is usually not synchronised with the movie from which the MAD visual features were extracted, mainly due to the varied durations of intro and outro of different movie source. Since we rely on the MAD visual features, the synchronisation is an essential step.

The pipeline for our automated data collection pipeline is briefly introduced in main paper Sect. 4. A schematic is shown in Fig. A.1, in detail:

1. We transcribe the mixed audio file using *WhisperX* [5] which provides accurate punctuated transcriptions with word-level timestamps.

2. The transcript is tokenized into sentences using the nltk python toolbox [11], resulting in transcription sentences and their corresponding temporal segments (inferred the start and end time of the first and last word in the sentence respectively).

3. Each sentence segment is assigned a single speaker identity (*e.g.* SPEAKER_00, SPEAKER_01, *etc.*) by performing speaker diarization on the mixed audio, whereby each sentence timestamp is provided as oracle voice activity detection. Specifically, we use SpeechBrain ECAPA-TDNN voice embeddings [24] trained on VoxCeleb [56] and Agglomerative Clustering with a threshold of 0.95.

4. To automatically identify the cluster associated with the AD speaker, we exploit the third-person nature of AD narrations and select the cluster with the lowest proportional occurrence of first- & second-person pronouns, *e.g.* "I" and "you" with 95 or more speaker segments.

5. To synchronise the segment timestamps with the original audio track from which the MAD visual features were extracted, we follow [79] and calculate the time delay $\tau$ between the original movie audio files and the mixed audio files via FFT cross-correlation. The timestamps of the identified AD segments are shifted according $\tau$ in order to synchronise them to the visual features and subtitles collected in MAD.

### A.2. AD Collection Pipeline for AudioVault

The collection pipeline for AudioVault is introduced in paper Sect. 5, we provide more details here. To collect text-only AD annotations from AudioVault, the final synchronisation step is unnecessary. Therefore, we follow steps 1-4 of the MAD denoising pipeline as described above, which takes as input the mixed audio tracks and outputs the ASR with timestamps from the possible AD speaker.

The large-scale collection from AudioVault audio files is noisy, *e.g.* some ADs are of lower-quality or are sourced from short movies. Therefore, we apply a stricter filtering step that removes movies containing fewer than 100 AD narrations or a word frequency of first- & second-person pronouns larger than 5%.

## A.3. Comparison with MAD-v1

The key advantages of our pipeline are three-fold: (1) it relies on *audio-based* speaker separation to identify the AD speaker among the movie characters, whereas the pipeline in the original MAD work [79] relies on *text-based* speaker separation by using the timestamps from the DVD subtitles and assume any ASR transcription outside of these timestamps is AD. The error is propagated because the official subtitles are non-exhaustive (some dialogue is missed by the official subtitles). (2) It requires only the mixed audio as input, whereas MAD must also source the official DVD subtitles and align them – presenting additional scaling costs and challenges. (3) It uses an advanced ASR model Whisper [66] which gives much more accurate transcriptions than the previous methods, especially on punctuation and the spelling of names and other identities.

## B. Qualitative Examples of MAD-v2 vs MAD-v1.

More qualitative examples of MAD-v2 and MAD-v1 is shown in Fig. A.2 and A.3. It is clear that our pipeline gives more accurate AD compared to the original MAD-v1, particularly in the spelling of names and the exclusion of dialogues.

**(a)**

**Manual Verification** With a dead-eyed stare, Chris sits in a cell.
**MAD-v1** <span style="color:red">Bring him back right.</span> with a dead eyed stare, Chris sits in a cell.
**MAD-v2 (ours)** With a dead-eyed stare, Chris sits in a cell.

**(b)**

**Manual Verification** Chris puts the rucksack on the floor.
**MAD-v1** Chris puts the <span style="color:red">rock psych</span> on the floor.
**MAD-v2 (ours)** Chris puts the rucksack on the floor.

**(c)**

**Manual Verification** Later he sits in a diner with Christopher.
**MAD-v1** Later he sits in a <span style="color:red"><?></span>.
**MAD-v2 (ours)** Later he sits in a diner with Christopher.

**(d)**

**Manual Verification** He comes up the steps.
**MAD-v1** He comes up <span style="color:red">at</span> steps. <span style="color:red">Can.</span>
**MAD-v2 (ours)** He comes up <span style="color:red">with</span> steps.

**(e)**

**Manual Verification** Chris looks ill as he watches Mr. Frohm's cab pull away.
**MAD-v1** Chris looks <span style="color:red">sailors</span> he watches Mr <span style="color:red">from</span> cab pull away.
**MAD-v2 (ours)** Chris looks ill as he watches Mr. <span style="color:red">From</span>'s cab pull away.

**(f)**

**Manual Verification** An uneasy look flickers across Chris' face as Jay leaves the washroom.
**MAD-v1** An uneasy look flickers across Chris's <span style="color:red">faces. Jail eats</span> the washroom.
**MAD-v2 (ours)** An uneasy look flickers across Chris' face as Jay leaves the washroom.

Figure A.2. **Comparison of the AD quality from MAD-v2 with MAD-v1.** The erroneous transcriptions are marked in <span style="color:red">red text</span>. 'Manual Verification' means we manually transcribe the AD narration from the audio track. The sample is originally from *The Pursuit of Happyness* (2006). The failure mode of MAD-v1 in each example is **(a)** dialogue leakage, **(b)** incorrect ASR, **(c)** missing words, **(d)** dialogue leakage, **(e)** incorrect ASR and name spelling, **(f)** incorrect name spelling.

15

**(a)**

**Manual Verification** Sully adjusts his seat harness.
**MAD-v1** Sully adjusts his seat harness, I.
**MAD-v2 (ours)** Sully adjusts his seat harness.

**(b)**

**Manual Verification** A male passenger looks up from his magazine.
**MAD-v1** Oh yeah, a male passenger, looks up from his magazine.
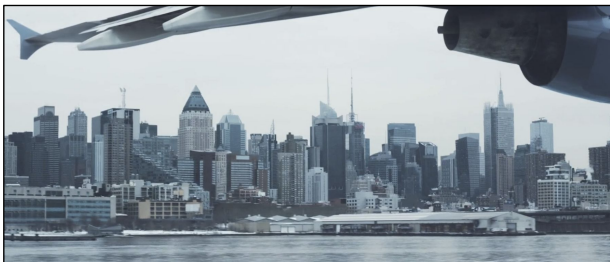**MAD-v2 (ours)** A male passenger looks up from his magazine.

**(c)**

**Manual Verification** Skiles turns to Sully in surprise.
**MAD-v1** The hudson skiles turns to sully and surprise i.
**MAD-v2 (ours)** Skiles turns to Sully in surprise.

**(d)**

**Manual Verification** A sightseeing helicopter comes into view over the dark waters of the Hudson.
**MAD-v1** <?> Sightseeing helicopter comes into view over <?>.
**MAD-v2 (ours)** A sightseeing helicopter comes into view over the dark waters of the Hudson.

**(e)**

**Manual Verification** The Manhattan skyline appears just under the wings of Flight 1549.
**MAD-v1** The manhattan skyline appears just under the wings of flight fifteen. Forty nine.
**MAD-v2 (ours)** The Manhattan skyline appears just under the wings of Flight 1549.

**(f)**

**Manual Verification** Sully sticks out an arm as the jet bellies down onto the river.
**MAD-v1** <?> The jet bellies down onto <?>.
**MAD-v2 (ours)** Sully sticks out an arm as the jet bellies down onto the river.

Figure A.3. **(continue) Comparison of the AD quality from MAD-v2 with MAD-v1.** The erroneous transcriptions are marked in red text. 'Manual Verification' means we manually transcribe the AD narration from the audio track. The sample is originally from *Sully: Miracle on the Hudson* (2016). The failure mode of MAD-v1 in each example is **(a)** dialogue leakage, **(b)** dialogue leakage, **(c)** dialogue leakage and incorrect ASR, **(d)** missing words, **(e)** number spelling and sentence partitioning, **(f)** missing words.

## C. Quantitative Comparison between MAD-v2 vs MAD-v1 on Grounding

We re-purpose the CLIP zero-shot video-language grounding (VLG) performance from [79] as an indicator of dataset quality. In detail, for both MAD-v2 and MAD-v1, we randomly choose a set of 5 movies from the *training split*, and compute the VLG performance with frozen CLIP visual and textual encoders. The AD textual quality and timestamps are the only factors that differ in this comparison. We use the MAD training split because we did not modify the val/test splits, which are from LSMDC annotations. The code to compute VLG performance is from `https://github.com/Soldelli/MAD`. The result in Table A.1 shows MAD-v2 annotations also benefit the VLG task.

| R@50 | IoU@0.1 | IoU@0.3 | IoU@0.5 |
|------|---------|---------|---------|
| MAD-v1-`Unnamed` | 32.08 | 22.85 | 14.26 |
| MAD-v2-`Unnamed` | **33.25** | **24.22** | **15.58** |

Table A.1. CLIP zero-shot VLG performance on MAD-v1 and MAD-v2.

## D. Additional Implementation Details

**Evaluation Metrics.** We use the `pycocoeval` package provided by `https://github.com/tylin/coco-caption` to compute the ROUGE-L, CIDEr, SPICE and METEOR. The package post-processes both the predicted text and ground-truch text internally to remove the punctuation and make them lowercased. To compute the BertScore, we use the package provided by `https://github.com/Tiiiger/bert_score`. Note that before computing the BertScore, both the predicted text and the ground-truth text are converted to lowercase without any punctuation, as these are factors that the BertScore is sensitive to.

**Alternative approach for vision-language fusion.** We investigate an alternative vision & language fusion mechanism whereby the context AD sentence prompts are fed as *language features* rather than *raw text tokens*. Empirically, we observe that raw text inputs outperform language features (e.g. 12.6 CIDEr in Table 2 vs. about 8.0 CIDEr when feeding language features).

## E. Additional Qualitative Examples

More qualitative examples are shown in Fig. A.4. It shows that our AutoAD model gives reasonable descriptions that fits movie description domain, like the actions (swim, dance), face expression (eyes widen). Note that under the oracle setting, our model is capable to learn the character names (sample **a, c, d, f**) mainly due to the extra information from the ground-truth context. Our model is still limited in its ability to identify characters accurately, *e.g.* in sample **f**, the movie shows Nick and *Daisy* are dancing. Whereas the oracle prediction describes that Nick and *Gatsby* are dancing, and the recurrent model simply predicts that a man and woman are dancing. Also the pronouns often appear in the recurrent prediction, such as the word 'his' in sample **a** and **b**, which shows the model learns the bias of pronouns but cannot recognize characters correctly.

## F. Dataset Splits

To clarify the dataset split of the movies in MAD and LSMDC, we list the movie IDs of each split we used (and not used) in this paper. The splits can also be found on our website `https://www.robots.ox.ac.uk/~vgg/research/autoad/`.

**MAD-v2.** It consists of 488 movies and all of them are used for training. We provide the cleaner ADs for these movies using the automated pipeline described above. This set is the same as the training set of movies proposed in MAD [79]. The movie IDs are:

```
[2723, 2730, 2731, 2735, 2738, 2745, 2750, 2758, 2768, 2778, 2787, 2800, 2801, 2814, 2818, 2854, 2869, 2870, 2873,
2911, 2913, 2928, 2934, 2944, 2948, 2970, 2986, 2992, 2996, 3001, 3014, 3020, 3021, 3023, 3033, 3040, 3049, 3050,
3059, 3060, 3066, 3070, 3103, 3106, 3113, 3114, 3117, 3129, 3138, 3146, 3153, 3160, 3170, 3171, 3209, 3239, 3253,
3276, 3277, 3295, 3314, 3339, 3340, 3354, 3376, 3393, 3401, 3408, 3414, 3417, 3447, 3464, 3480, 3482, 3500, 3509,
3510, 3513, 3521, 3548, 3575, 3590, 3599, 3611, 3625, 3720, 3743, 3759, 3773, 3820, 3834, 3837, 3858, 3905, 3911,
3922, 3977, 4001, 4007, 4010, 4017, 4031, 4043, 4053, 4061, 4071, 4080, 4082, 4143, 4156, 4200, 4204, 4210, 4253,
4266, 4299, 4303, 4305, 4368, 4377, 4378, 4390, 4423, 4434, 4451, 4455, 4460, 4480, 4489, 4528, 4535, 4551, 4576,
4578, 4587, 4596, 4597, 4608, 4611, 4618, 4634, 4635, 4638, 4644, 4664, 4670, 4671, 4684, 4702, 4709, 4719, 4728,
4740, 4741, 4753, 4772, 4778, 4797, 4798, 4813, 4815, 4839, 4880, 4884, 4888, 4901, 4902, 4914, 4925, 4929, 4933,
4936, 4950, 4962, 4970, 4977, 4982, 4992, 5014, 5041, 5055, 5063, 5074, 5093, 5101, 5118, 5139, 5144, 5217, 5236,
5237, 5257, 5259, 5265, 5270, 5283, 5293, 5308, 5335, 5366, 5367, 5369, 5417, 5420, 5432, 5449, 5461, 5469, 5473,
5477, 5494, 5506, 5510, 5511, 5522, 5563, 5565, 5568, 5574, 5575, 5577, 5583, 5594, 5605, 5607, 5634, 5641, 5649,
```

**(a)** **(b)** **(c)**
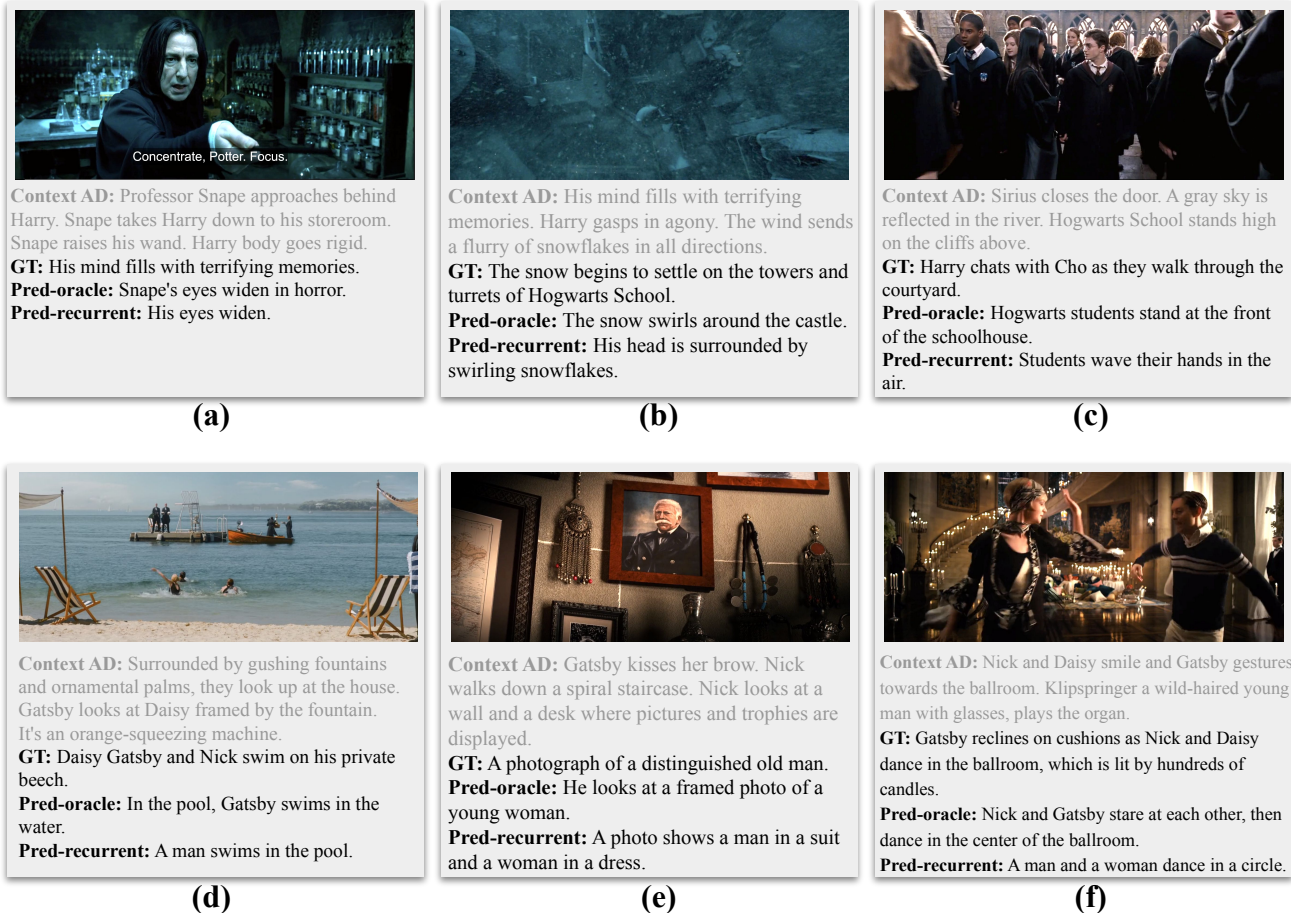
**(d)** **(e)** **(f)**

Figure A.4. **Qualitative examples of AutoAD model.** We show the ground-truth AD and the AD predictions under both the oracle and recurrent settings. Previous AD context is shown in gray. Samples are taken from *Harry Potter and the Order of the Phoenix* (2007) and *The Great Gatsby* (2013).

5677, 5678, 5682, 5685, 5700, 5735, 5737, 5743, 5749, 5752, 5758, 5762, 5792, 5807, 5814, 5818, 5819, 5828, 5852,
5865, 5872, 5873, 5898, 5900, 5913, 5923, 5950, 5958, 6012, 6013, 6022, 6048, 6055, 6057, 6076, 6086, 6090, 6137,
6153, 6154, 6156, 6177, 6186, 6194, 6224, 6232, 6319, 6334, 6394, 6402, 6491, 6521, 6607, 6613, 6617, 6629, 6636,
6655, 6656, 6672, 6685, 6701, 6706, 6741, 6769, 6770, 6775, 6810, 6811, 6816, 6819, 6832, 6833, 6837, 6859, 6869,
6870, 6878, 6890, 6952, 6959, 6992, 6994, 7001, 7005, 7007, 7026, 7036, 7050, 7055, 7131, 7195, 7196, 7243, 7682,
7882, 8152, 8276, 8295, 8346, 8496, 8578, 8587, 8589, 8593, 8598, 8601, 8608, 8616, 8618, 8637, 8734, 8766, 8767,
8811, 9110, 9277, 9380, 9384, 9386, 9387, 9419, 9421, 9451, 9456, 9460, 9461, 9462, 9481, 9482, 9488, 9502, 9504,
9509, 9510, 9515, 9519, 9526, 9528, 9529, 9535, 9552, 9555, 9575, 9576, 9583, 9595, 9606, 9615, 9617, 9618, 9619,
9620, 9638, 9642, 9644, 9647, 9654, 9659, 9676, 9689, 9719, 9724, 9732, 9733, 9735, 9737, 9738, 9741, 9747, 9750,
9751, 9754, 9756, 9761, 9773, 9774, 9785, 9799, 9846, 9896, 9906, 9920, 9952, 10142, 10149, 10202, 10322, 10527,
10536, 10784, 10813, 10836, 10861, 10894, 10965, 11003, 11010, 11099, 11129, 11139, 11140, 11143, 11147, 11148, 11154,
11318, 11321, 11345, 11396, 11430, 11438, 11530, 11620, 11727, 11796, 11962, 12010, 12079, 12090, 12125, 12131, 12132,
12144, 12147, 12148, 12186, 12211, 12220, 12222, 12263, 12273, 12294, 12324, 12358, 12504, 12563, 12585, 12618, 12653,
12658, 12743, 12852, 12869, 12900, 12906, 12911, 12923, 12958, 13018, 13027, 13031, 13045, 13140, 13146, 13159, 13165,
13187, 13191, 13201]

**MAD-eval Evaluation Set.** It consists of 10 movies, which are obtained by $set$(MAD val/test) $\cap$ $set$(LSMDC val), that excluding LSMDC train movies for the ease of future comparison. The annotations are inherited from the LSMDC dataset, and we use both the *named* and *unnamed* version of it, where the *named* version can be downloaded from the LSMDC website[2]. The movie IDs are:

---
[2]https://sites.google.com/site/describingmovies/download?authuser=0

[1005_Signs, 1026_Legion, 1027_Les_Miserables, 1051_Harry_Potter_and_the_goblet_of_fire, 3009_BATTLE_LOS_ANGELES, 3015_CHARLIE_ST_CLOUD, 3031_HANSEL_GRETEL_WITCH_HUNTERS, 3032_HOW_DO_YOU_KNOW, 3034_IDES_OF_MARCH, 3074_THE_ROOMMATE]

**Unused MAD movies.**   151 movies from MAD val/test are *not used* in either training or testing in our paper. They are the intersection $set$(MAD val/test) $\cap$ $set$(LSMDC train/test). The movie IDs are:

[0001_American_Beauty, 0002_As_Good_As_It_Gets, 0003_CASABLANCA, 0004_Charade, 0005_Chinatown, 0006_Clerks, 0007_DIE_NACHT_DES_JAEGERS, 0008_Fargo, 0009_Forrest_Gump, 0010_Frau_Ohne_Gewissen, 0011_Gandhi, 0012_Get_Shorty, 0013_Halloween, 0014_Ist_das_Leben_nicht_schoen, 0016_O_Brother_Where_Art_Thou, 0017_Pianist, 0019_Pulp_Fiction, 0020_Raising_Arizona, 0021_Rear_Window, 0022_Reservoir_Dogs, 0023_THE_BUTTERFLY_EFFECT, 0026_The_Big_Fish, 0027_The_Big_Lebowski, 0028_The_Crying_Game, 0029_The_Graduate, 0030_The_Hustler, 0031_The_Lost_Weekend, 0032_The_Princess_Bride, 0033_Amadeus, 0038_Psycho, 0041_The_Sixth_Sense, 0043_Thelma_and_Luise, 0046_Chasing_Amy, 0049_Hannah_and_her_sisters, 0050_Indiana_Jones_and_the_last_crusade, 0051_Men_in_black, 0053_Rendezvous_mit_Joe_Black, 1001_Flight, 1002_Harry_Potter_and_the_Half-Blood_Prince, 1003_How_to_Lose_Friends_and_Alienate_People, 1004_Juno, 1006_Slumdog_Millionaire, 1007_Spider-Man1, 1008_Spider-Man2, 1009_Spider-Man3, 1010_TITANIC, 1011_The_Help, 1012_Unbreakable, 1014_2012, 1015_27_Dresses, 1017_Bad_Santa, 1018_Body_Of_Lies, 1019_Confessions_Of_A_Shopaholic, 1020_Crazy_Stupid_Love, 1028_No_Reservations, 1031_Quantum_of_Solace, 1033_Sherlock_Holmes_A_Game_of_Shadows, 1034_Super_8, 1035_The_Adjustment_Bureau, 1037_The_Curious_Case_Of_Benjamin_Button, 1038_The_Great_Gatsby, 1039_The_Queen, 1040_The_Ugly_Truth, 1042_Up_In_The_Air, 1043_Vantage_Point, 1045_An_education, 1046_Australia, 1047_Defiance, 1048_Gran_Torino, 1050_Harry_Potter_and_the_deathly_hallows_Disk_One, 1052_Harry_Potter_and_the_order_of_phoenix, 1054_Harry_Potter_and_the_prisoner_of_azkaban, 1055_Marley_and_me, 1057_Seven_pounds, 1058_The_Damned_united, 1059_The_devil_wears_prada, 1060_Yes_man, 1061_Harry_Potter_and_the_deathly_hallows_Disk_Two, 1062_Day_the_Earth_stood_still, 3001_21_JUMP_STREET, 3002_30_MINUTES_OR_LESS, 3003_40_YEAR_OLD_VIRGIN, 3004_500_DAYS_OF_SUMMER, 3005_ABRAHAM_LINCOLN_VAMPIRE_HUNTER, 3007_A_THOUSAND_WORDS, 3008_BAD_TEACHER, 3012_BRUNO, 3013_BURLESQUE, 3014_CAPTAIN_AMERICA, 3016_CHASING_MAVERICKS, 3017_CHRONICLE, 3018_CINDERELLA_MAN, 3020_DEAR_JOHN, 3022_DINNER_FOR_SCHMUCKS, 3023_DISTRICT_9, 3024_EASY_A, 3025_FLIGHT, 3026_FRIENDS_WITH_BENEFITS, 3028_GHOST_RIDER_SPIRIT_OF_VENGEANCE, 3030_GROWN_UPS, 3033_HUGO, 3035_INSIDE_MAN, 3036_IN_TIME, 3037_IRON_MAN2, 3038_ITS_COMPLICATED, 3039_JACK_AND_JILL, 3040_JULIE_AND_JULIA, 3041_JUST_GO_WITH_IT, 3042_KARATE_KID, 3043_KATY_PERRY_PART_OF_ME, 3045_LAND_OF_THE_LOST, 3046_LARRY_CROWNE, 3047_LIFE_OF_PI, 3048_LITTLE_FOCKERS, 3049_MORNING_GLORY, 3050_MR_POPPERS_PENGUINS, 3051_NANNY_MCPHEE_RETURNS, 3052_NO_STRINGS_ATTACHED, 3053_PARENTAL_GUIDANCE, 3054_PERCY_JACKSON_LIGHTENING_THIEF, 3055_PROMETHEUS, 3056_PUBLIC_ENEMIES, 3058_RUBY_SPARKS, 3060_SANCTUM, 3061_SNOW_FLOWER, 3062_SORCERERS_APPRENTICE, 3063_SOUL_SURFER, 3066_THE_ADVENTURES_OF_TINTIN, 3067_THE_ART_OF_GETTING_BY, 3069_THE_BOUNTY_HUNTER, 3070_THE_CALL, 3071_THE_DESCENDANTS, 3072_THE_GIRL_WITH_THE_DRAGON_TATTOO, 3073_THE_GUILT_TRIP, 3075_THE_SITTER, 3076_THE_SOCIAL_NETWORK, 3077_THE_VOW, 3078_THE_WATCH, 3079_THINK_LIKE_A_MAN, 3081_THOR, 3082_TITANIC1, 3083_TITANIC2, 3084_TOOTH_FAIRY, 3085_TRUE_GRIT, 3086_UGLY_TRUTH, 3087_WE_BOUGHT_A_ZOO, 3088_WHATS_YOUR_NUMBER, 3089_XMEN_FIRST_CLASS, 3090_YOUNG_ADULT, 3091_ZOMBIELAND, 3092_ZOOKEEPER]