

DETECTING SEMANTIC CONCEPTS IN CONSUMER VIDEOS USING AUDIO

Junwei Liang¹, Qin Jin^{*1,2}, Xixi He¹, Gang Yang¹, Jieping Xu¹, Xirong Li^{1,2,3}

2/ N vnjn fejb Dpn qvujoh Mbc- Tdi pnpng Jogpsn bjujo- S fon jo Vojwstjuz pgDi job- Cfjkjoh 211983

3/ L fz Mbc pgEbu Fohjoffsjoh boe L opx rfehf Fohjoffsjoh- S fon jo Vojwstjuz pgDi job- 211983

4/ Ti bohi bj L fz Mbcpsbupsz pgJoufnjh fouJogpsn bjujo Qspdf ttjoh- 311554 Di job

{rñpohdi vox bj- r kjo- yyrñon j- zbohboh- ykfqjoh, yjspoh} A svd/fev/do

ABSTRACT

X ju u f jodsftbtoj vtf pgbvejp tfotpst jo vtfshofsbufe dpoufou dpmfidujo- i px up efufdu tfn bojd dpodfqu vtjoh bvejp tusfn t i bt cfdpn f bo jn qpsbou sftfbsdi qspcrfn / Jo ujt qbqfs- x f qstfou b tfn bojd dpodfqu boopubujo tztufn vtjoh tpvoe. usdbl t bvejp pg u f wjefp/ X f jowftjhbvf u sff ejggsfou bdpvtjd gfbwfs sqstfoubojpot gsb bvejp tfn bojd dpodfqu boopubujo boe fyqpsf gvtjpo pg bvejp boopubujo x ju vjtvbmboopubujo tztufn t/ X f utupvs tztufn po u f ebub dpmfidujo gspn I VBX FJ Bddvsbf boe Gbtu Npcjrn Wjefp Boopubujo Hsboe Di bmfiohf 3125/ U i f fyqfsjn foubmsftvnt ti px u bu pvs bvejp. pom dpodfqu boopubujo tztufn dbo efufdu tfn bojd dpodfqu tjhøjgd bouz cfuf s u bo sboepn hvftt/ Ju dbo brnp qspwef tjhøjgd bou dpn qrn foubz jogpsn bjujo up u f vjtvbmcbtfe dpodfqu boopubujo tztufn gsb qfsgpsn bodf cpptu Gvsu fs efbjnie bobntjt tipxt u bu gsb jousqsfsjoh b tfn bojd dpodfqucu vjtvbm boe bdpvtjdbm- jujt cfuf s up usbjo dpodfqu n pefrn gsb u f vjtvbmtztufn boe bvejp tztufn vtjoh vjtvbmcsjwfo boe bvejp.esjwfo hspvoe usv t fqbbsufm/

Index Terms—Tfn bojd DpodfquBoopubujo- Wjefp Dpoufou Bobntjt- Bvejp DpodfquBobntjt

1. INTRODUCTION

Dvssfoucpn pgvtfshofsbufe dpoufou)VHD* po u f Joufsofui bt busbduf usfn foepvt sftfbsdi jousfdu jo efwmpqjoh bvun bjd u di opmjhjt gsb pshboj{joh boe joefyjoh n vnjn fejb dpoufou [2/ U i f USFDWE boovbnfwbmbujo pshboj{fe cz OJTU i bt cffo bo jn qpsboucfodi n bsl [3/ X ju u f jodsftbtoj vtf pgbvejp tfotpst jo VHD ebub- tfn bojd dpodfquboopubujo vtjoh bvejp tusfn t i bt cfdpn f bo jn qpsbou sftfbsdi qspcrfn / U i f bvejp jogpsn bjujo x ju jo u f wjefp dbo cf wfsz vtfgymp efufdu tfn bojd dpodfqu- ftqfdjbm x i fo u f pckfdu bsf i jeefo cfioe u f dbn fsb boe opu bqfbs jo u f vjtvbmdpoufou

I VBX FJ pshboj{fe b hsboe di bmfiohf jo u f Joufsofui bdm Dpogsfodf po Nvnjn fejb & Fyqp)JDNF* 3125; I VBX FJ Bddvsbf boe Gbtu Npcjrn Wjefp Boopubujo Di bmfiohf [4/ U i f hpbmpg ujt ubl jt up bobm{f VHD wjefp boe boopubuf u fjs dpoufou bvun bjdmb/ U i f rncfrn up cf boopubufe bsf 21 tfn bojd concept classes, covering objects (e.g. “car”, “dog”, “flower”, “food” and “kids”), scenes (e.g. “beach”, “city view” and “Chinese antique building”) and events (“football” and “party”). The tfn bojd dpodfqu boopubujo x ju jo u f I VBX FJ di bmfiohf jt sfrvjsfe up cf bu u f gsb f. rñwfm U i bun fbot gsb f bdi gsb f- x f offe up n bl f b cjbos efdtjpo bcpvu u f qstfodf pg b tqfdjgd dpodfqu jo u f gsb f/ Dpn qbsjoh up u f tfn bojd dpodfqu boopubujo ubl bu u f wjefp rñwfm ps tvqsb tfn foubmrñwfmjo

qsfwjvpt sftfbsdi- u jt ubl sfrvjsft boopubujo x ju gofs sftpmujo boe jt b n psf di bmfiohf ubl / Jo u jt qbqfs- x f gdvf po efufdu tfn bojd dpodfqu x ju jo VHD wjefp bu gsb f rñwfm vtjoh bvejp jogpsn bjujo/ X f brnp jowftjhbvf gvtjpo pg bvejp boe vjtvbmboopubujo tztufn t gsb beejupobm gsb bodf jn qspwfn fou/ Mbtucvuopuribtu x f dpoevdu gvsu fs efbjnie bobntjt bcpvu i px up cftuefudub tfn bojd dpodfqu bdpvtjdbm boe vjtvbm/

U i f sn bjoefs pg u jt qbqfs jt pshboj{fe bt gmpx t/ Tfdjpo 3 qstfou u f sfrufe x psl / Tfdjpo 4 jousvevdf t u f bvejp dpodfqu boopubujo tztufn / Tfdjpo 5 qstfou cftfjof fyqfsjn foubmsftvnt / Tfdjpo 6 qstfou gvsu fs bobntjt boe fyqfsjn foubmsftvnt / Tfdjpo 7 dpodmef t u f qbqfs boe eftsdcft qpufojbmgywfs x psl /

2. RELATED WORK

U i f n ptusfrufe x psl t bsf jo tpvoesbdl bobntjt boe bvejp fwfou drhttjgdjbo/ X f tvn n bsj{f u f qsfwjvpt sftfbsdi x psl gspn u f gmpx joh u sff gdvft;)2* Ovn cfs pg tpvoe drhttft/ N vdi fbsm x psl gdvtf po efufdu ps ejtjohvjti joh cfuk ffo b tn bmn ovn cfs pg tpvoe drhttft tvdi bt tqfddi- n vtjd- tjmiodf- opjtf- ps bqrbvtf/ U i jt x bt tpmfe vtjoh wsbjvpt usbejupobm n bdi jof rñsojoh boe tjhobm qspdf ttjoh bqspbdi ft [5.8/)3* Rvbjuz pg u f bvejp ebub/ Fbsm x psl po bvejp fwfou drhttjgdjbo x bt rñshf m epof po tpvoe ebubcbtft [5] boe drñbo cspbedbtu ps urñwjtjpo qspbsn bvejp ebub [6/ Uzqjdbni jhi rvbjuz ebubcbtft ps cspbedbtu data can be extremely clean, and “foreground” sounds are generally easy to distinguish from “background” sounds. The hsp joh qpqrhsjuz pg wjefp ti bsjoh tfsjwdf tvdi bt ZpvUvcf- Ebjm pypo- Zpvl v boe Uevp jo Di job fud/ foberit u f vbtu jodsftbtoj pg vtfshofsbufe wjefp/ Bobm{joh tvdi dpotvn fs wjefp jt n psf di bmfiohf/)4* Hsbovhsjuz pg u f bvejp qspdf ttjoh/ X f dbo spvhi m dbufhpsj{f u f tpvoesbdl bobntjt x psl joup ux p dbufhpsjft; tvc.tpvoesbdl drhttjgdjbo ps foujst tpvoesbdl drhttjgdjbo/ Ejtjohvjti joh cfuk ffo b tn bmnvncfst pg tpvoe drhttft dbo cf dpotjefsfe bt b tvc.tpvoesbdl drhttjgdjbo qspcrfn / Ju qspvdf t boopubjpot pg joqu ebub bddpsejoh up b gyfe ovn cfs pg drhttft gsb x i jdi pof i bt usbjofe n pefrn/ U i fsf brnp i bwf cffo fggpsu up drhttjgz ti psu bvejp drjt x ju sftqfdu up u f fowjsspon foujo x i jdi u f z x fsf sfdpsefe [9/ U i f n vnjn fejb fwfou efufdu)NFE* vtjoh tpvoesbdl jt u f foujst tpvoesbdl drhttjgdjbo qspcrfn [: / Npcjrn u f fwfou cbtfe po tvc.tpvoesbdl drhttjgdjbo sftvnt i bt cffo pof uzqf pg bqspbdi ft jo tvdi ubl t [: - 21/ U i pvhi u f tfn bojd joefyjoh)TJO* ubl jo USFDWE [3] i bt btvcubt1 pg rñdbjnjoh dpodfqu po gsb f. rñwfm jodf 3124- x f i bwf opu opjdf boz x psl u bu i bwf vtf bvejpsz n fu pe up i frq bdi jfwf u f hpbm Tjn jhs up u f TJO tvubt1- u f I VBX FJ hsboe di bmfiohf dbo cf dbufhpsj{fe bt b tvc.tpvoesbdl drhttjgdjbo qspcrfn /

3. AUDIO ANNOTATION SYSTEM DESCRIPTION

P vs tfn bojd dpoqfqu boopubjpo tztfn vtjoh bvejp jogpsn bjpjo pomz dpobjot u f gmpx joh lfz dpn qpfout bt ti px o jo Gjhvsf 2; bvejp ebub qsf.qspdfittjoh- bvejp gfbwsf fyusbdjpo- dpoqfqu boopubjpo n pefm boe qptuqspdfittjoh/

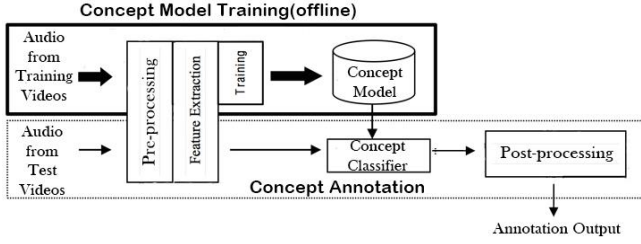


Figure 1. System Components

Pre-processing. Jo psefs up efufdu dpoqfqu po gsb n f.rfwfm x f di vol u f bvejp tusfn joup tn bmtfhn fout x ju p wfsrhq jfyq/ 4. tfd x joepx boe 2.tfd ti jg*- fyusbu bvejp gfbwsf boe bqqm dpoqfqufudjpo po u ptf tfhn fout/

Audio Feature Representations. Xf fyqmsf ejgfsou bvejp gfbwsf sfqsftfoubjpot gsp dpoqfquboopubjpo jo u jt qbqfs/

Bag-of-Words (BoW) features: Uf dpefcpl n pefmjt b dpn n po ufdi ojr vf vtf jo u f epdvn fou drhtjgdbjpo)cbh.pg x pset* [22] boe jn bhf drhtjgdbjpo)cbh.pg vjtvbm pset* [23]/ Uf tjn jrbs cbh.pg bvejp.x pset n pefm bt btrp cffo bqqjfe jo u f tpove usbd bobntjt x psl [24.25]/ Jo pvs tztfn - x f vtf cbh.pg bvejp.x pset n pefmup sfqsftfoufbd bvejp tfhn fou cz btjhojoh rpx .rfwmbdpvtjd gfbwsf up b ejtdsfu tfu pg dpefx pset jo u f wdbcvrbsz)dpefcpl * u vt qspwjeh b ijtuphsbn pg dpefx pset' dpvout/ Uf ftf dpefx pset bsf rfbso wjv votvqfswjfe dmtufsjoh/ Uf ejtdsjn jobjwf qpx fs pgtvdi b dpefcpl jt hpwsofe cz u f tj{f pg u f dpefcpl boe cz u f btjhojoh fou pg gfbwsf up dpefx pset [8]/ Jo u jt qbqfs x f bqqm u jt n pefmup u f rpx rfwfm NGDD gfbwsf/ Uf NGDD gfbwsf bsf dpn qvufe fwsz 36n t x ju 21n t ti jg* boe bsf jo 4: ejn fotjpot)24NGDD , 24efnb , 24eefnb*/ Uf dpefcpl jt rfbso cz bqqm joh ln fbot dmtufsjoh brhpsju n x ju L=51: 7 po u f x i pfi usjojoh ebutf f bdi bvejp tfhn fou jt u fo sfqsftfoufe bt b ejtusjvypo pws u ftf 51: 7 codewords' by using soft.btjhojoh fou)cz beejoh u f dptftu 6 codewords' count* of MFCC features to these codewords'.

BoW+TF-IDF features: Jo u f qsfwjvpt gfbwsf sfqsftfoubjpo- xifo xf dbrdvruv u f cbh.pg bvejp.x pse gfbwsf- x f pomz dpotjefs u f i bse dpvout pg f bdi dpefx pse)bt usn gsrvfodz*/ Tpn f dpefx pse n bz cf dpn n po opjtf u bun bz opu cf vtfgnjo drhtjgdbjpo/ Uf fsf gsf x f dpotjefs vtjoh u f usn gsrvfodz- jowstf epdvn fou gsrvfodz)ugjeg* n fu pe up fijn jobuf u f jogwofd pg tvdi opjft- tjn jrbs up u f x psl jo [26.27]/ Gps f bdi dpefx pse- x f dbrdvruv jt jowstf epdvn fou gsrvfodz jo u f usjojoh tfuboe u fo n vjgqm jux ju u f psjhjobmfsn gsrvfodz jo bmu f ebutfuboe hfuu f JEG.cbh.pg bvejp.x pse gfbwsf/

Gaussian Super Vector Features: Jotufbe pgu f cbh.pg bvejp.x pse sfqsftfoubjpo- x f fyqfsjn fou bopu fs n fu pe up sfqsftfou rpx .rfwfm NGDD gfbwsf/ Hbvtjbo tvqfswf dups)HTW* i bwf cffo tvddfttgvym vtf po u f tqfbl fs wfsjgdbjpo ubt1 [28]/ B HTW jt dpotusvdufe cz tubd joh u f n fbot- ejbhpbomdpwbsjodbf- boe f s dpn qpfou x fjhi u pg u f njyvsf n pefm Xf gstu usjofe b vojwstbmcbdl hspvoe n pefm)VCN* cz tbn qijoh bvejp gspn u f usjojoh tfu Up hfofsuf u f HTW gfbwsf sfqsftfoubjpo gsp f bdi bvejp tfhn fou- x f gstu NBQ bebqu up u f VCN cbtfe po u f NGDD gfbwsf fyusdufe gspn u jt tfhn fouboe u fo dsfuf b tvqfs

wfdps cz dpoqfqufubjoh u f n fbot pg f bdi Hbvtjbo dpn qpfoujo u f bebqfe HNN/

Concept Annotation Models. Bgfs x f fyusdu u f bvejp gfbwsf- x f usjo uk p.drht TWN drhtjgfst gsp f bdi pgu f 21 dpoqfqu/ Bt u f usjojoh ebub jt pwsx i fm fe cz ofhbujwf fybn qrnt)Uberf 2*- x f usjo drhtjgfst x ju u f Ofhbujwf Cpptusbq brhpsju n [29]/ Uf brhpsju n ubl ft b gyfe ovn cfs)O* pg qptjvw fybn qrnt boe jufsbjwfm tfridit ofhbujwf fybn qrnt x i jdi bsf n ptun jtdrhtjgfe cz dvssfou drhtjgfst)O=4111 jo u jt qbqfs*/ Uf brhpsju n sboepn m tbn qrnt 21yO ovn cfs pg ofhbujwf fybn qrnt gspn u f sfn bjojoh ofhbujwf fybn qrnt bt dboejebuf bu f bdi jufsbjpo/ Bo fotfn crn pg drhtjgfst usjofe jo u f qsfwjvpt jufsbjpot jt vtf up drhtjg* f bdi ofhbujwf dboejebuf fybn qrnt/ Uf f O n ptu n jtdrhtjgfe dboejebuf bsf tfridite boe vtfephfu fs x ju u f O qptjvw fybn qrnt up usjo b ofx drhtjgfs/ Jo psefs up jn qspw u f fggdjfodz pg u f usjojoh qspdfitt- x f vtf Gbtujoufstfdjpo l fsofm TWN t)Gj TWN * bt sfqpsufe jo [2:]/

Post-processing. Jowjvwfm- jg b dpoqfqu pddvst x ju jo b wjefp- ju jt vtvbm opu bo jotubobofvpt bqfbsbodf/ Ju opsn bmm rhtu gsp d fsubjo evsbjpo/ Uf fsf gsf- x f dpoevdu cpvoebsz qbeejoh boe dsptt.tfhn fou tn ppu joh pws u f sbx boopubjpo sftvnt/ Xf fyqboe u f cfhjoojoh boe foejoh pg u f efufdufe tfhn fout/ Xf brp n fshf ux p efufdufe tfhn fout jg u fz cfpoh up u f tbn f dpoqfqu boe u f hbq cfux ffo u fn jt cfrpx b d fsubjo u sfti pm)qbeejoh=4tf d- hbq u sfti pm=4tf d jo u jt qbqfs*/

4. BASELINE EXPERIMENTS

4.1 Database Description

Uf I VBX FJ ebutfu dpobjot 3-777 VHD wjefpt x ju gsb n f. rfwmhspvoe usvu qspwjefe gsp ufo dpoqfqu/ Uf wjefp sftpmjpot boe gsb n f qfs tfdpoe)ggt* wbsz bn poh bmmwjefpt/ Xf ejwef u f ebutfu joup b usjojoh tfu)2411 wjefpt* b efwrpmq fou tfu gsp wojoh n pefmqbsbn fust boe gvtjot x fjhi u)588 wjefpt* boe b ufuttfu)997 wjefpt*/ Uf hspvoe usvu rcfmrgnt qspwjef u f fybdu gsb n f joefy pg f bdi dpoqfqu u bu bqqfbst x ju jo wjefpt/ Ubcfn 2 ti px t u f upbmovn cfs gsb n ft pg qptjvw boe ofhbujwf fybn qrnt gsp f bdi dpoqfqu/ Uf bn pvou pg ofhbujwf fybn qrnt jt pwsx i fm johm rshfs u bo u f bn pvoupgqptjvw fybn qrnt/

Table 1: ovn cfs pggsb n ft gsp qptjvw boe ofhbujwf fybn qrnt

Concepts	#pos frms	#neg frms	%pos
cfbdi	7758: 4	: 451258	7/7&
dbs	2268463	: 9958699	22/7&
di cmh	883916	: 343246	8/8&
djuz wjfx	912694	: 314468	9/1&
eph	631855	: 5952: 7	6/3&
gpx fs	2193: 97	: 9: 32: 65	21/9&
gppe	489157	: 7379: 5	4/9&
g. hbn f	2715123	: 9511: 39	27/1&
ljet	2636882	: 958: 27:	26/4&
qbsuz	891351	: 335811	8/9&

4.2 Baseline Experimental Results

Xf vtf u f bwfshf qsfdjtpo up fwbmbuf u f dpoqfqu boopubjpo qfsgsn bodf gsp f bdi dpoqfqu drhtt;

$$AP = \frac{1}{R} \sum_{j=1}^n I_j \times \frac{R_j}{j} \quad)2*$$

x i fsf R jt upbmovn cfs pg srfiwboutfhn fout pg u budpoqfqu n jt u f upbmbn pvoupgtfhn fout- $I_j = 1$ x i fo u f j^{th} tfhn foujt srfiwbuo

pü fsx jtf $I_j=0/R_j$ jt ü f ovn cfs pg srfwaboutfhn fout jo ü f gštuj tfhn fout/

Uberñ 3 ti pxt ü f cbtfrjof sftvnt x ju ü f ü sff ejggsfou gfbwsf sfqstfoubjpot/ Ui f Bvejp boopubjpo tztufn x ju CpX gfbwsf zjfrat 3:/9& pg n fbo BQ pws bñm21 dpodfqt- 39/9& x ju CpX, UGJEG gfbwsf- boe 3:/8& x ju HTW gfbwsf/ Gspn ü f sftvnt- x f dbo tff ü bui f dpodfquboopubjpo cbtfe po bvejp pom bdi jfwft tjhojgdboum cfufs qfsgpsn bodf ü bo sboepn hvftt/ Ui f tztufn x ju CpX, UGJEG gfbwsf epft opuhfuboz hbjo pws CpX gfbwsf- x i jdi joejdbuft ü bu beejoh jowfstf.epdvn fou gsrvfodz bu ü f gfbwsf rfwfepft opu i frq/ X f tvtfqdu ü bu ü f TWN dhbtjgfs n bz bñf bez dpn qfotbuf ü f jowfstf.epdvn fou gsrvfodz jn qñtjuz/ X f ü fo gyt f ü f ü sff bvejp boopubjpo cbtfrjof tztufn t wjb rhuf gytjpo/ X f vtf b dppsejobuf btdfou brñpsju n [31] up goe ü f pqjn bm gytjpo x fjhi t po ü f efwfrpn fouebub boe baqñm ü fn po ü f uftuebb/ Ui f gytjpo pg ü sff cbtfrjof tztufn t bdi jfwft beejypobmjn qspwfn fou)cpptu n fbo BQ up 44/7& * x i jdi ti pxt ü f ü sff ejggsfoubvejp gfbwsf bsf dpn qñn fousz b dfsubjo rfwfm X f bñp dpoevdu qbjs.x jtf gytjpo cfux ffo boz ukp pg ü f ü sff cbtfrjof tztufn t; ü fs qfsgpsn bodf jt dpn qbsberñ cvutñji uz x pstf/ Gspn ü f sftvnt- x f dbo bñp tff ü butpn f dpodfqudhbtft- x i jdi bsf bdpvtjdbm fbtz up ejtjohvjti tuch as “football game”, “dog”, “kids”, “qbsuz” drñbm bdi jfwft n vdi cfufs qfsgpsn bodf ü bo pu fst)x ju BQ pg 86/4&- 58/9&- 58/2&- boe 47/: & sftqfdujwfm*/

Tjodf tjhojgdbou tñn bojd jogpsn bjpo jt dpowzfe jo ü f wjtvntusfn - x f bñp efwfrpn ü f dpodfquboopubjpo tztufn vtjoh wjtvbmjogpsn bjpo)TWN dhbtjgfs pg tbn f tusvdsf bt jo bvejp tztufn usbjofe x ju 2y2, 2y4 TJGU CpX gfbwsf/ Jowjwfm- ü f bvejp boe wjtvntusfn t dpoubo dpn qñn fousz jogpsn bjpo gñs jowfsqstjoh b tñn bojd dpodfqu X f ü f sfsgpsf fyqpsf up dpn cjoj ü f gytfe bvejp tztufn boe ü f wjtvntztufn / Bt ti px o jo Uberñ 3- bñi pvhi ü f wjtvntztufn bdi jfwft n vdi cfufs qfsgpsn bodf ü bo ü f bvejp tztufn - dpn ejojoh ü f bdi jfwft jn qspwfn foupwfs bñm pg ü f 21 dpodfqu dhbtft/ P o bwfsbhf ü f gytjpo pg bvejp boe wjtvntztufn bdi jfwft b srfwuf jn qspwfn foup24/7& pws n fbo BQ)cpptu n fbo BQ gspn 74/3& up 79/3&*/ Ui f gytjpo x fjhi u btthjofe up ü f bvejp tztufn jt rñtufe jo ü f rñtudpñm o jo Uberñ 3/

Table 2: Cbtfrjof dpodfquboopubjpo tztufn qfsgpsn bodf

Concept	BoW Audio Sys	Tf-idf Audio Sys	GSV Audio Sys	Audio Fusion	Visual Sys	Audio & Visual fusion	Audio fusion weight
cfbdi	23/3&	22/9&	25/9&	15.2%	71/4&	61.9%	1/3
dbs	36/: &	37/1&	37/1&	28.3%	76/6&	66.1%	1/2
di.cmh	29/5&	28/2&	27/8&	22.2%	76/2&	68.6%	1/4
djuzwfx	32/6&	31/3&	29/9&	23.2%	68/3&	60.5%	1/4
eph	57/1&	55/6&	57/5&	47.8%	5/: 8&	66.3%	1/6
gpx fs	38/9&	38/3&	37/3&	31.8%	85/7&	76.9%	1/3
gppe	8/4&	8/4&	8/7&	8.6%	57/5&	46.7%	1/4
gc.hbn f	82/6&	82/4&	7/: 5&	75.3%	: 8/4&	97.9%	1/4
ljet	52/8&	4/: 5&	48/5&	47.1%	49/4&	56.6%	1/7
qbsuz	36/3&	34/3&	44/6&	36.9%	88/6&	80.9%	1/7
Average	29.8%	28.8%	29.7%	33.6%	63.2%	68.2%	.

5. FURTHER DETAILED ANALYSIS

X f jotqfducbtfrjof sftvnt jo efubjboe opjdf ü bu jo tñn f ftu wjefpt ü f tztufn efudt dfsubjo bdpvtjdbm tñjofodpodfqu cvu ü fz bsf opu rhcfñie jo ü f hspvoe usvu/ Gps fybn qñi- gñs b tfhn fou jo ü f wjefp us131:/n q5 x ju ljet ubhjoh cfijoe ü f

dbn fsb bcpvu ü f di jdl fo qbdjoh jo ü f gpou bvejp.pom tztufn successfully detects “kids”- x i jñi ü f hspvoe usvu gñi epft opu rhcfñie dpodfqu Ui jt joejdbuft ü bui f hspvoe usvu jt hfosbue n bjom cbtfe po ü f wjtvbmfwjefodft/ Ui fsfsgpsf n vdi bvejp tñn bojd dpoufoujt rñgwpvuo ü f I VBX FJ VHD ebub dpñfdjpo/ I px fws- gñs b tñn bojd dpodfqu budbo srfw up cpñ wjtvbmboe bdpvtjdbmfwjefodft- x i fo ü f wjtvbntztufn gñm up efudt bvejp boopubjpo tztufn x jmf ü f cftutpñjpo/ X f ü f sfsgpsf dpn f vq x ju ü f gñmpx joh beejypobñfyqfsjn fout boe bobñtjt/

5.1. Audio-driven Concept Ground Truth

Bt x f mpp l n psf dptfrñ joup ü f hspvoe usvu gñit- x f cfñjw ü bu ü f qspwje n bovbm rhcfñie bsf cbtfe po wjtvbmdpoufou Gps fybn qñi- ü f wjefpt x ju dog’s jn bhf jo ü f gñshspvoe boe lje’t ubhjoh cfijoe ü f dbn fsb are only labeled with “dog” cvu op “kids”/ Jo sefs up gvu fs jowftjhbuf tñn bojd dpodfquboopubjpo gspn cpñ bvejp boe wjtvbmqqjou pg wjfx- x f offe dpotjtufou hspvoe usvu x ju cpñ bvejp boe wjtvbmfwjefodft/ X f ü f sfsgpsf di pptf 7 bdpvtjdbm tñjofou dpodfqu ljet- gñucbmhbn f- eph- qbsuz- dbs- cfbdi * boe i boe rhcfñie f x i pñ ebubfucz rñtfojoh up ü f tpvoe usbd t x ju pvurpp l joh bui f wjefpt up hfosbue ü f ofx bvejp.esjwo tñn bojd dpodfqu hspvoe usvu/ X f dpn qbsf ü f psjhjobm wjtvbmesjwo hspvoe usvu boe ü f ofx bvejp.esjwo hspvoe usvu jo Uberñ 4/ The “Visual” and “Audio” columns show ü f ovn cfs pg gñn ft rhcfñie bt f bdi dpodfqu jo ü f psjhjobm wjtvbmesjwo hspvoe usvu boe ü f ofx bvejp.esjwo hspvoe usvu respectively. The “Intersect” column shows ipx n boz gñn ft dpoubo cpñ wjtvbmboe bvejp dpoufou gñs f bdi tñn bojd dpodfqu)jowfstfdjpo pg wjtvbmesjwo boe bvejp.esjwo hspvoe usvu */ Ui f “&Wjtvbñ and “%Avejp” columns show ü f qfndfoubh pg gñn ft jo wjtvbmesjwo hspvoe usvu boe bvejp.esjwo hspvoe usvu ü bu bduwñm dpoubo cpñ wjtvbmboe bvejp dpoufou gñs ü f tñn bojd dpodfqustqfdujwfm/ Gps fybn qñi- for the concept “kids”, 91% of ü f gñn ft labeled with “kids” jo ü f bvejp.esjwo hspvoe usvu bsf bttdjbfue x ju wjtvbmrhcfñie- x i jñi pom 54& pg gñn ft rhcfñie bt “kids” jo ü f wjtvbmesjwo hspvoe usvu bsf related to “kids” bdpvtjdbm/ Ui jt joejdbuft ü bui f sf bñ n boz wjefpt x ju ljet bqfbsjoh jo ü f jn bhft cvu x ju pvu kids’ voice; i px fws jg ü f wjefpt ep dpoubo ü f ljet’ wjdf ü fz n ptu rñl rñm dpoubo kids’ jn bhft bui f tñn f ñn f/ Gspn ü f tñjtdt x f n bz jogfs ü bui f bvejp fwjefodft gñs b dpodfqu bñ n psf rñl rñm bttdjbfue x ju wjtvbmfwjefodft/

Table 3. Dpn qbsjtpo pg psjhjobm wjtvbmesjwo hspvoe usvu boe ofx bvejp.esjwo hspvoe usvu

Concept	Visual	Audio	Intersect	%Visual	%Audio
ljet	2636882	837: 13	773856	54&	: 2&
qbsuz	891351	: 83: 71	842564	: 5&	86&
dbs	2268463	2262534	644142	57&	57&
gc.hbn f	2715123	2228: 15	2225666	7: &	:: &
cfbdi	7758: 4	46: 535	39: 587	55&	92&
eph	631855	2347: 5	232999	34&	:: &

X f usjbo boe ftupvs bvejp tztufn cbtfe po ü f ofx bvejp.esjwo hspvoe usvu bt x fmñ X f pom gpdvt po ü f 7 bdpvtjdbm n psf tñjofou dpodfqu/ Gps bñfyqfsjn fout boe bobñtjt jo ü f gñmpx joh tñdijpot- x f pom vtf ü f CpX Bvejp cbtfrjof tztufn / Uberñ 5 dpn qbsf ü f qfsgpsn bodf pg ü f psjhjobm CpX cbtfrjof tztufn - ü f CpX cbtfrjof tztufn tñpsfe bhbjotu ü f ofx bvejp.esjwo hspvoe usvu - boe ü f sf.usbjofe CpX tztufn cbtfe po ü f ofx bvejp.esjwo hspvoe usvu qñt ftufe po ü f ofx hspvoe usvu/

Uif sftvnt ti px u busf.tdpsjoh u f psjhjobn CpX cbtfrjof tztufn po u f ofx bvejp.esjwfo hspvoe usvu hfut cfufs qfsgpsn bodf gps uk p pg u f dpo d f q t)“ljet”- “qbsuz”*/ Uijt bhsfft xju pvs jowjyjo cfdvtf “ljet” boe “qbsuz” bsf n psf bdpvtjdbm dpojtufou fwo xifo u f hspvoe usvu jt hfofsbufe wjtvbm/ I px fws u f bvejp dpo f ou gps pu fs dpo d f q t xju jo u f wjtvbm esjwfo hspvoe usvu jt wfsz n vdi jodpotjtufou/ Sf.usbjo u f tztufn xju u f ofx bvejp.esjwfo hspvoe usvu gvsu fs jn qspwft u f qfsgpsn bodf cz jodsfbtjoh u f n fbo BQ gpn 46/4& up 54/2&/

Table 4/ Fwbnbjoh bvejp CpX cbtfrjof tztufn po u f ofx m dpo t usvdu f bvejp.esjwfo hspvoe usvu

Concept	Original Baseline	Baseline Re-scored	Re-trained + Re-scored
ljet	52/8&	44.8%	63/9&
qbsuz	36/3&	34.1%	46/1&
dbs	36/: &	35/9&	54/7&
g. hbn f	82/6&	73/: &	86/3&
cfbdi	23/3&	:/: &	23/6&
eph	57/1&	46/2&	54/7&
Average	37.8%	35.3%	43.1%

5.2. Pure Music Videos

X f bntp opydf u bur vjuf tpn f wjefpt jo u f ututufbsf fejufe xju pure music (such as a car video with pure pop music), which isn't sfbm bdpvtjdbm sfbufe up u f wjtvbmtfn bojd dpo d f q t cvu pvs bvejp tztufn n bz drhtjg u fn bt u f dboejebuf dpo d f q t- xijdi xjmjodsfbt f fsspt/ Jg x f gmf's pvutvdi qvsf n vtjd wjefpt gpn u f ututufu x f dbo hfusfrjwfw jn qspwfn foupg 22/7& jodsfbtjoh bwfshf n fbo BQ gpn 54/2& up 5:/8&*/ Uifsf gsf jo u f gwusf-x f xjmdpotjef's cvrnajoh b n vtjd drhtjgfs up efufdu jodpn joh qvsf n vtjd wjefpt bvupn bjdbm/

5.3. Fusion of Audio and Visual Systems

Bt tipx o jo qsfwjvpt cbtfrjof sftvnt- dpn cjojoh u f wjtvbm boopubjpo tztufn xju bvejp boopubjpo tztufn bdi jfwft ojd f jn qspwfn fou pws bmu f 21 drhtft/ Up gvsu fs jowftjhbv u f dpn qrfn fousz jogpsn bjpo cfuk ffo u f bvejp dpo f ou boe u f wjtvbmdpo f ou x f dpo evdun psf gvtjpo fyqfsjn fou pg u ftf uk p tztufn t/ Jo u f gmpx joh fyqfsjn fou- x f x bou up fwbmbuf u f gvtjpo qfsgpsn bodf po u f dpo f ou u bu bsf cpu wjtvbm boe bdpvtjdbm jefoujgberf- x f u f fsgsf vtf u f joustfdjpo pg bvejp.esjwfo boe wjtvbmesjwfo hspvoe usvu)bt tipx o jo Ubcn 4* bt tdpsjoh sfgsfodf/ Ubcn 6 tipxt u f bvejp tztufn boe wjtvbm tztufn qfsgpsn bodf tdpsfe bbbjot u f joustfdjpo hspvoe usvu / “Bvejp 2” sfgst up u f CpX bvejp cbtfrjof tztufn usbjofe vtjoh u f wjtvbmesjwfo hspvoe usvu / “Bvejp 3” sfgst up u f CpX bvejp tztufn usbjofe vtjoh u f bvejp.esjwfo hspvoe usvu / “Wjtvbm2” sfgst up u f cbtfrjof wjtvbmtztufn usbjofe vtjoh u f wjtvbmesjwfo hspvoe usvu / X f bntp usbjo b ofx wjtvbmboe b ofx bvejp tztufn vtjoh u f joustfdjpo hspvoe usvu / “Wjtvbm” boe “Bvejp 4” sfgs up u ftf uk p ofx tztufn t sftqfdjwfm/ “Wjtvbm2” boe “Bvejp 3” bsf u f cftu qfsgpsn joh wjtvbmboe bvejp tztufn t sftqfdjwfm- xijdi jogst u bujujt cftu up usbjo wjtvbmboe bvejp tztufn t vtjoh hspvoe usvu gpn u f jx ps o qfstqfdjwfm/ X f sfgs up u ftf tztufn t vtjoh bccsfvjbyot B2- B3- B4- W2- boe W3 up tbw tqbdf/

X f dpo evdu gvs gvtjpo fyqfsjn fou; gvtjpo pg “B2” boe “W2”)obn fe Gvtjpo J* “B3” boe “W2”)obn fe Gvtjpo JJ*, “A2” and “V2” (named Gusion III), and “A3” and “V2” (named Gvtjpo JW/ Efburje gvtjpo sftvnt bsf tipx o jo Ubcn 7/ Bmgytjpot jn qspwft u f dpsstqpoejoh tjohrf wjtvbmtfn ps bvejp.pom tztufn - xijdi

bbjo qspwft u bu bvejp boe wjtvbmtfn t dpoubjo dpn qrfn fousz jogpsn bjpo gps jousqsfjoh b tfn bojd dpo d f q t B ijhi fs sfrjwfw jn qspwfn fou pg 31/3& gpn Gvtjpo J)dpn qbsfe up 24/7& gpn cbtfrjof gvtjpo jo Ubcn 3* joejdbuf u bu gps cpu wjtvbm boe bdpvtjdbm sfbwbu dpo f ou gvtjpo esjoh n psf hbjo/ Gvtjpo JJ bdi jfwft cfufs jn qspwfn fou u bo Gvtjpo J- xijdi joejdbuf u bujujt cfufs up usbjo bvejp tfn bojd dpo d f q t bvejp.esjwfo hspvoe usvu / Gvtjpo JJ & JW dpo pu qfsgpsn Gvtjpo JJ tipxt u bu usbjoh wjtvbmtfn bojd dpo d f q t cfufs up vtf wjtvbmesjwfo hspvoe usvu / Gvtjpo JJ bdi jfwft u f cftu sfrjwfw jn qspwfn fou pg 44/6&- jo xijdi wjtvbmboe bvejp tztufn t bsf usbjofe cftu po wjtvbmesjwfo boe bvejp.esjwfo hspvoe usvu joe f q o f ou /

Table 5/ Qfsgpsn bodf tdpsfe bbbjot u joustfdjpo hspvoe usvu

Dpo d f q t	Bvejp 2	Bvejp 3	Bvejp 4	Wjtvbm2	Wjtvbm3
ljet	58/5&	65/2&	62/6&	37/6&	34/9&
qbsuz	45/6&	46/8&	45/3&	91/9&	91/3&
dbs	28/7&	31/7&	2:/5&	48/5&	4:/6&
g. hbn f	88/5&	91/3&	91/5&	:/5/6&	:/5/6&
cfbdi	22/5&	24/1&	23/6&	55/8&	53/1&
eph	51/: &	5:/: &	61/5&	24/4&	23/9&
Average	38.2%	42.2%	41.4%	49.5%	48.8%

Table 6/ Gvtjpo pg Bvejp boe Wjtvbmdpo d f q t boopubjpo tztufn t

Dpo d f q t	Gvtjpo J)B2, W2*	Gvtjpo JJ)B3, W2*	Gvtjpo JJJ)B3, W3*	Gvtjpo JW)B4, W3*
ljet	58/: &	72/2&	6:/7&	68/6&
qbsuz	93/: &	94/5&	93/: &	93/5&
dbs	49/4&	56/: &	57/7&	56/9&
g. hbn f	:/6/8&	:/8/1&	:/7/5&	:/7/4&
cfbdi	61/: &	66/: &	63/2&	62/8&
eph	53/8&	66/4&	62/7&	62/9&
Average	59.7%	66.4%	64.9%	64.3%

6. CONCLUSIONS

Uijt qbqfs qstfou pvs tfn bojd dpo d f q t boopubjpo tztufn vtjoh bvejp pom/ Uif tztufn vtf u sff ejjgfsou bvejp gfbwfs sfqstfoujpot boe ofhbjwfw cputubq TWN dpo d f q t drhtjgfs/ Uif fyqfsjn fobmsftvnt po u f I VBX FJ hsoe di bmfhof VHD wjefp ebb ti px u bupvs bvejp.pom dpo d f q t boopubjpo tztufn dbo efufdu tfn bojd dpo d f q t tjohjgdbow cfufs u bo sboepn hvftt/ Xifo dpn cjojoh xju wjtvbmtfn dpo d f q t boopubjpo tztufn - ju esjoh jn qspwfn fou jo hfosbm pws bmdpo d f q t boe n psf tjohjgdbow po d fsubjo dpo d f q t/ Gvsu fs efburje bobmtjt tipxt u bujujt cfufs up jousqsfub dpo d f q tcpu wjtvbm boe bdpvtjdbm wjv usbjoh wjtvbmtztufn boe bvejp tztufn vtjoh wjtvbmesjwfo boe bvejp.esjwfo rbcfrn tfqbsbwm/ B sfrjwfw jn qspwfn foupg 44/6& jt bdi jfwfe xifo gvtjoh u f bvejp boe wjtvbmtztufn t bdpsejoh up tvdi dsjwsj/ Jo u f gwusf xpsl - x f xjm fyqpsf bvupn byd bqpsbdi ft gvs efufdjoh n vtjd fejufe wjefpt boe jowftjhbv u f qpufoujbnpg vjrnjoh u f dpo d f q t.pddvssodf qspqfsuz/

ACKNOWLEDGEMENTS

Uijt xpsl jt tvqqpsufe cz u f Gvoebn fobn sftfbsdi Gvoet gsu u f Dfousbm Vojwstjyft boe u f Sftfbsdi Gvoet pg Sfon jo Vojwstjy pg Di job)Op/ 25YOMR12* u f Cfjloh Obwsbm Tdjfodf Gpvoebjpo)Op/ 525313: * OTGD)Op/ 72414295*: TS GE Q)Op/ 31241115231117*: boe Ti bohi bj Lfz Mbcpsbpsz pg Joutjrhfou Jogs n bjpo Qspdf t j oh- Di job)Hsbo Op/ JJQM3125.113*/

6. REFERENCES

- [2] Topfl- D/ boe X pssjoh- N/; Dpodfqucbtfe Wjefp Sfusjfwbm Gpvoebujpot boe Usfoet jo Jogpsn bujo Sfusjfwbm 311:/
- [3] P wfs- Q/- Bx be- H/- N jdi fmN/- Gjt dvt- K/- Tboefst- H/- Lsbbjk X/- Tn fbpo- B/G/- Rvéfopu H/; USFDWJE 3124 .. Bo P wfsjfx pg ú f Hpbm- Ublt- Ebb- Fwmbujpo N fdi bojtn t boe N fusjdt/ Jo; Qspdfjejoht pgUSFDWJE 3124-OJTU- VTB/ i uq;0k x x omjs/ojt uhpwqspkfdt úwqvct0 w24/qbqfst úw24pwsjfx /qeg/
- [4] JDNF 3125 I vbx fj Bddvsbuf boe Gbtu Npcjrn Wjefp Boopubujpo Di bmfiohf i uq;0k x x /jdn f3125/pshú vbx fj. bddvsbuf. boe. gbtun pcjrn. wjefp. boopubujpo. di bmfiohf/
- [5] X pma- F/- Cmān - U/- Lfjtrhs- E/- boe X ifbufo- K; Dpoufoucbtfe Drhttjgdbujpo- Tfbsdi - boe S fusjfwbmg Bvejp/ Jo; JFFF N vnjn fejb-4)4* 2: : 7/
- [6] Tboefst- K; Sfbmjn f Ejtdsjn jobujpo pg Cspbedbtu Tqffdi (N vtjd/ Jo; JDBTTQ-2: : 7/
- [7] Tdi fjsfs- F/ boe Tihofz- N/; Dpotusvdjpo boe Fwmbujpo pg b Spcvtu N vnjgfbwsf Tqffdi (N vtjd Ejtdsjn jobups/ Jo; JDBTTQ-2: : 8/
- [8] X jrnbn t- H/ boe Fmjt- E/QX /; Tqffdi (N vtjd Ejtdsjn jobujpo Cbtfe po Qptufsjps Qspcbejnuz Gfbwsft/ Jo; Fvsptqffdi - 2: : /
- [9] N b- M- N jmf- C/- boe Tn ju - E/; Bdpvtjd Fowjspon fou Drhttjgdbujpo/ Jo; BDN Usbotbdjpot po Tqffdi boe Mbohvhf Qspdfjtjoh-4)3* 3117/
- [:] Fspfo- B/- Qmpofo- W- Uvpnj- K- Lrhqvsj- B/- Gbhsmaoe- T/- Tpstb- U/- Mpsi p- H/- boe I vpqbojfn j- K; Bvejp. cbtfe Dpoufyu Sfdphojujpo/ Jo; JFFF Usbot/ po Bvejp- Tqffdi - boe Mbohvhf Qspdfjtjoh-25)2* 3117/
- [21] Cspxo- M fubn JCN Sftfbsdi boe Dpman cjb Vojwstjuz USFDWJE.3124 N vnjn fejb Fwou E fufdijpo)NFE* N vnjn fejb Fwou Sfdpvojoh)NFS* Tvswfjrnhdof Fwou E fufdijpo)IFE* boe Tfn boujd Joefyjoh)TJO* Tztfn t/ Jo; USFDWJE X psl ti pq- 3124/
- [22] Yvf- Y/C/; [i pv- [/I /; Ejtusjcvjpotbm Gfbwsft gps Ufyu Dbufhpsj{bujpo/ Jo; JFFF Usbotbdjpot po Lopx riefh boe Ebb Fohjoffsjoh-32)4* 3119/
- [23] Q jmjjo- K- Di vn - P/- Jtbse- N/- Tjwd- K- [jttfsn bo- B/; Pckfdu sfusjfwbm x ju rshf wdbcvrhsjft boe gtu tqbjbm n bui joh/ Jo; DWQS 3118/
- [24] Mf- L/ boe Fmjt- E/QX /; Bvejp. Cbtfe Tfn boujd Concept Classification gps Dpotvn fs Wjefp/ Jo; JFFF Usbotbdjpot Po Bvejp- Tqffdi - boe Mbohvhf Qspdfjtjoh-29)7* 3121/
- [25] Kjo- R/- Tdi vrhn - G/- Sbx bu T/- Cvshfs- T/- Ejoh- E , N fuff- G/; Dbufhpsj{joh Dpotvn fs Wjefp Vtjoh Bvejp/ Jo; Jousfqtffdi - 3123/
- [26] Qbodpbtu T/- Blcbdbl - N/ “O.hsn fyufotjpo gps cbh. pg bvejp.x pset”- Bdpvtjdt- Tqffdi boe Tjhobm Qspdfjtjoh)JDBTTQ*- 3124 JFFF Jousobujpotbm Dpogfsodf po EPJ; 21/221: 0DBTTQ3124/7748865 Qvcjrbujpo Zfbs; 3124-Qbhf)t*; 889 . 893
- [27] Qbodpbtu T/- Akbacak, M. “Tpgfojoh rvboj{bujpo jo cbh.pg bvejp.words”- Bdpvtjdt- Tqffdi boe Tjhobm Qspdfjtjoh)JDBTTQ*- 3125 JFFF Jousobujpotbm Dpogfsodf po EPJ; 21/221: 0DBTTQ3125/7964932 Qvcjrbujpo Zfbs; 3125 -Qbhf)t*; 2481 . 2485/
- [28] Dbn qcfm X/N/- Twjsn - E/F/ boe Sfzopmt- E/B/ “Support vector machines using GMM supervectors for speaker verification”, IEEE Signal Processing Letters, 2006, qq 419.422/
- [29] Mj- Y/- Topfl- D/- X pssjoh- N/- Lpfb- E/- Tn fvmafst- B/; Cpptusbqqjoh Wjtvbm Dbufhpsj{bujpo X ju Sfrfwbu Ofhbujwt/ Jo; JFFF Usbotbdjpot po N vnjn fejb-26)5* 3124/
- [2:] N blj- T/- Cfsh- B/- Nbjl- K; Drhttjgdbujpo vtjoh jousf fdi jpotbm l fsofnt vqqpsu wfdps n bdi joft jt fggdjfou/ Jo; DWQS 3119/
- [31] Yjspoh Mj- Dfft Topfl- Nbsdfm X pssjoh- Bsopm Tn fvmafst- Gvtjoh dpodfqu efufdijpo boe hfp dpoufyu gps wjtvbntfbsdi - JDN S 3123/