

用于高真实感口型合成的唇区肌肉模型

周 维 汪增福

(中国科学技术大学自动化系,合肥 230027)

摘 要 为了快速合成真实感强的口型,在 Waters 肌肉模型的基础上,提出了一种新的唇区肌肉模型,弥补了 Waters 模型对于唇区肌肉的较复杂运动描述不完全的不足。该模型以面部解剖学为基础,通过面部运动的解剖学机理的研究,将唇区的整体运动表示为若干个子运动的线性组合。该模型可用于语音驱动的动画合成。此时,只需在说话人的唇区标定少数几个特征点,就可以获得一组唇区肌肉参数,进而建立相应的肌肉模型。借助于该模型,可以在唇区附近的线性肌的联合作用下,合成说话时的各种口型。实验结果表明,该肌肉模型不仅计算代价小,且合成的口型真实感较强,具有很强的实用性。

关键词 Waters 肌肉模型 唇区肌肉模型 子运动 语音动画 线性肌

中图法分类号:TP391.9 文献标识码:A 文章编号:1006-8961(2008)11-2238-06

Muscle-based Lips Model for Realistic Lips Synthesis

ZHOU Wei, WANG Zeng-fu

(Department of Automation, University of Science & Technology of China, Hefei 230027)

Abstract Based on the Waters' muscle model(WMM), a novel lip muscle model(LMM) is proposed in this paper. LMM perfects the inaccurate description of the complicated lip muscle movements in WMM. On the basis of facial anatomy, the global lip movement is divided into a few sub-movements. The reconstruction of the lip movement is based on the linear combination of the sub-movements. In the application of modeling talking face, several feature points are marked to obtain a group of lip parameters. All kinds of lip shapes are synthesized using the proposed LMM and the adjacent linear muscle model. The experimental results show that the proposed model is practical in view of its low computational cost and ability of producing all kinds of realistic synthesized lip shapes.

Keywords Waters' muscle model, lip muscle model, sub-movement, talking face, linear muscle

1 引 言

3D 人脸建模和面部运动合成一直是计算机图形学领域中的重要和热点问题,同时具有真实感的人脸动画合成在人机交互、电影和游戏制作、医学辅助研究、视频会议等领域中扮演着重要角色。自 20 世纪 70 年代 Parke 等人所做的开创性工作^[1]以来,人脸动画合成技术获得了长足的发展,人脸运动建

模的方法日益多样化。其中,已有的建模方法大体可分为基于表演驱动的方法^[2]、基于物理学运动模型的方法^[3]、基于肌肉参数模型的方法^[4]和基于学习的方法^[5]等 4 类。

其中,基于表演驱动的方法是运用专门的硬件设备将人脸的变形转换成合成的人脸模型,并运用由专门设备所探测到的人脸运动参数来控制合成头部模型的运动,但这种方法对硬件设备要求很高,代价相对昂贵;基于物理学运动模型的方法是将动态

基金项目:模式识别国家重点实验室开放基金和中国科学技术大学与中国科学院自动化研究所智能科学与技术联合实验室开放基金项目(JL0602)

收稿日期:2006-11-01; **改回日期:**2007-04-23

第一作者简介:周 维(1981 ~),男。2007 年获中国科学技术大学博士学位。主要从事人脸语音动画合成研究与数字图像处理研究。

E-mail: zhouwei8@mail.ustc.edu.cn

质点-弹簧系统用于人脸动画的合成,首先建立人脸的质点网格结构,然后通过控制质点-弹簧系统的运动产生面部动画。虽然依据这种模型所生成的人脸动画真实感很强,但计算量相对较大;基于学习的方法不考虑脸部的运动机理,而是运用主成分分析(PCA)技术直接从一组数字化人脸数据出发来得到所期望的合成人脸,该方法虽可以合成具有极高真实感的人脸动画,但由于学习效果与所使用的数据库的质量密不可分,因此需要的代价也是相当昂贵的。

Waters 在 1987 年提出了一种基于解剖学的肌肉模型^[4],用于表现质点在肌肉力的作用下产生的位移。该模型先将面部肌肉分为 3 类,然后分别对它们进行建模,并通过各类肌肉的组合运动来构造 FACS (facial action coding system) 中的 AU (action unit) 单元,进而合成面部运动。Waters 提出的肌肉模型从生理解剖学角度出发,对面部肌肉从运动本质上进行数学描述。该肌肉模型有以下特点:①肌肉模型的建立不需要表演驱动建模中专门的硬件仪器支持;②与物理学运动模型相比,其计算代价更小,且易于动画的实时合成;③相比于基于学习的方法,它不需要建立庞大的高质量图像数据库,且在参数计算上也相对简单,便于实时处理。综上所述,该肌肉模型由于它的简洁性和所构建的面部网格各单元之间所具有的独立性而被广泛使用。然而,Waters 提出的肌肉模型在某些情况下也有比较大的局限性。例如,在该模型中,轮匝肌被描述为一个具有径向收缩功能的肌肉单元,但在语音动画设计等实际应用中,该描述远远不能满足口型合成的需要。众所周知,语音动画是以人脸动画合成为基础的,不论是文本驱动^[6]、语音驱动^[7]或者是混合驱动^[8]的方法,大多是基于上述 4 种方法合成人脸动画,其真实感直接取决于人脸动画合成的真实感,而口型合成则是其中最为关键的部分。由于嘴唇运动往往比较复杂,使用上述简单的轮匝肌模型不能达到较精确描述说话人口型的目的,从而也不能取得逼真地刻画唇区运动的动画效果。在上述研究背景下,本文在 Waters 提出的模型的基础上,对轮匝肌及其附近的肌肉群进行了较为精确的建模,并依据所构建的肌肉模型,以较小的代价实现了各种口型较高自然度的合成。

2 人脸唇区解剖学

唇区建模涉及唇区中的骨骼、肌肉和皮肤的运动学建模。建模效果的好坏在很大程度上依赖于对人脸唇区中的上述组织在解剖学意义上的了解。

唇区的骨骼运动,主要指下颌骨的运动。下颌骨的运动包括开合运动、前后运动和侧方运动^[9]。所谓开合运动是指下颌的升降运动,其表现为绕两耳根处关节腔连线轴的旋转运动;前后运动主要是指下颌的前伸和后收运动;而侧方运动则主要是指在咀嚼食物时的小范围左右运动。人说话时,唇区的骨骼运动基本上表现为开合运动。

唇区的运动,其在物理上由相关肌肉组织所产生的组合运动来完成。每一个实际运动通常由若干个子运动组合而成。唇区主要肌肉组织的分布如图 1 所示。为了实现对上述组合运动的精确建模,本文采用对各个子运动分别进行建模的方法来建立唇区肌肉模型。从发音学角度进行归纳和整理后的子运动如下^[10]:

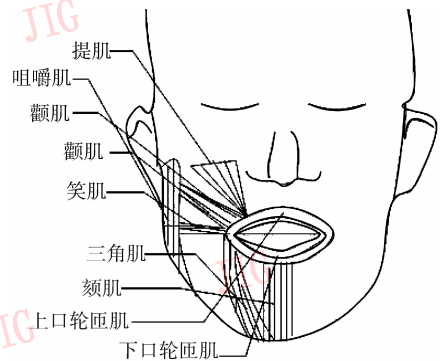


图 1 唇区的肌肉组织示意图

Fig. 1 Lip muscle tissue

(1) 闭唇和圆唇运动。与其相关的肌肉是口轮匝肌 (orbicularis oris)。它在闭唇运动中联合附近的其他肌肉,内收嘴唇,并使上嘴唇下移,下嘴唇上移,将嘴唇拉近牙齿。例如,发双唇辅音/b/。在圆唇运动中,通过径向收缩口轮匝肌,即可以将嘴唇变圆。

(2) 提升上唇运动:与其相关的肌肉是提肌 (levator labii superioris) 和颧肌 (zygomaticus)。运动时,提肌用于提升上嘴唇。例如,发唇齿摩擦音/f/。

(3) 降下唇运动:与其相关的肌肉是三角肌 (depressor anguli oris) 等,运动时,向侧下方拉下嘴

唇。例如,发元音/i/。

(4) 突唇运动:与其相关的肌肉是颞肌(mentalis)和口轮匝肌。运动时,颞肌将下巴上的皮肤往上拉,同时突出下唇,而口轮匝肌则使整个唇型变圆,并且可以翘起与突出上下唇。颞肌配合口轮匝肌运动,帮助突出和变圆嘴唇。例如,用于圆唇音/u/和汉语拼音中/zhi/等发音。

(5) 收缩嘴角运动:与其相关的肌肉是颊肌(buccinator)、笑肌(risorius)、颧肌和三角肌。颊肌收缩时牵引口角向侧后,使嘴唇贴近牙齿。一般用于唇齿音,还有/i/、/e/等元音的发音。笑肌运动时侧拉嘴角,并在发/i/、/e/等音时,用于辅助颊肌和颧肌向后拉嘴角。颧肌向侧上方拉嘴角,与提肌一起向上运动,用来发唇齿摩擦音,侧面运动用来发/s/音。三角肌运动时降嘴角,使嘴巴不要完全闭上,一般与颞肌一起运动。

3 唇区肌肉模型

Waters 提出的肌肉模型^[4]对人的面部肌肉进行了简化处理。由于这个原因,使得该模型虽然在一些基本表情的合成方面有较好的表现,但在较复杂的口型动作表现方面则存在比较大的局限性。例如,在 Waters 提出的肌肉模型中,口轮匝肌是用可以进行简单径向收缩的括约肌来表示,这种过分的简化却影响了对于较复杂的口型动作的合成效果。本文在对与发音相关的唇区解剖学进行深入研究的基础上,对口轮匝肌进行了较为精确的建模,并在 Waters 定义的其余线性肌肉的配合下,使之能够较为理想地合成各种说话中的复杂口型。

前已述及,Waters 用径向收缩的括约肌来表示口轮匝肌,并将其用于一些基本表情的合成。但是,人在说话时,口轮匝肌的作用不仅仅是产生径向收缩运动,还参与诸如突唇、翘唇等主动复杂运动的合成。不仅如此,它在和其附近的肌肉群的共同作用下,还会产生一些诸如下拉、侧拉等被动运动。为了在视觉上表现上述主动和被动运动,必须对口轮匝肌模型进行修正和改进。

口轮匝肌的整体运动可用它的各部分所产生的子运动的线性组合来表达,即

$$D_{Clo} = \alpha_{Rad} D_{Rad} + \alpha_{Clo} D_{Clo} + \alpha_{Ptr} D_{Ptr} + \alpha_{Til} D_{Til} + \alpha_{Add} D_{Add} + \alpha_{Oth} D_{Oth} \quad (1)$$

其中, D_{Clo} 为整体运动; D_{Rad} 为唇部径向收缩运动,是

口轮匝肌的基本运动,其在圆唇、突唇和翘唇运动中都会伴随发生; D_{Clo} 为闭唇运动; D_{Ptr} 为突唇运动; D_{Til} 为翘唇运动,其可看作是突唇运动的一种; D_{Add} 为内收嘴唇运动,其可看作是闭唇运动的一种; D_{Oth} 为由口轮匝肌附近的线性肌肉所产生的运动。 $A = \{\alpha_{Rad}, \alpha_{Clo}, \alpha_{Ptr}, \alpha_{Til}, \alpha_{Add}, \alpha_{Oth}\}$ 为参数集,各分量分别对应相关运动的权值。

这里,根据各个子运动的内在运动机理和外在表现形式,分别对其进行数学建模。建模过程如下:

首先,通过标定的特征点确定口轮匝肌的影响区域。如图 2 所示,按照 Moubaraki 等人的做法^[11],利用 3 个标定的特征点,就可以对唇区进行建模。其中, $M_{1,0}$ (上唇顶)、 $M_{2,0}$ (左唇角)和 $M_{3,0}$ (下唇底)3 个关键特征点可通过标定得到,而嘴唇上的其余控制点则由这 3 个点通过插值得到。如图 3 所示,建模的具体做法如下:首先指定口轮匝肌区域中的 3 个关键控制点 P_0 、 P_1 、 P_2 ,然后根据口轮匝肌区域的对称性得到嘴唇中心点 P_3 和右唇角点 P_4 。

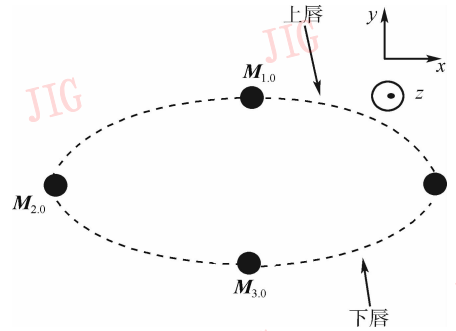


图 2 3 个关键特征点嘴唇模型

Fig. 2 The three key feather points of lip model

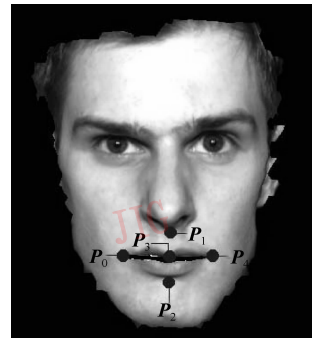


图 3 口轮匝肌区域特征点示意图

Fig. 3 The feather points of orbicularis oris

根据上面得到的关键控制点,就可以分别建立前文归纳的各子运动模型。其中,唇部径向收缩子模型和口轮匝肌附近的线性肌影响子模型沿用

Waters 提出的模型。

本文首先对闭唇运动进行建模。闭唇运动(例如发/b/、/p/、/m/音)时,轮匝肌径向收缩,上下唇前端唇缘区域分别以两唇区域的上下边界曲线(如图4椭圆中上下边界曲线)为轴,微向口内旋转。同时,上唇区域以上唇边界曲线为参考,微向下运动;下唇区域以下唇边界曲线为参考,微向上运动。因此,在闭唇运动中,上唇区域中的某一点 P 的运动由下式给出:

$$D_{\text{Clo}}(\mathbf{P}) = \mathbf{R}_{\text{Clo}}(\theta_p) + T_{\text{Clo}}(d_p)$$

$$\begin{cases} \theta_p \propto \frac{\|x_p x_{p_0}\| \cdot \|y_p y_{p_1}\|}{\|x_p x_{p_4}\| \cdot \|y_p y_{p_1}\|} & \text{如果 } \mathbf{P} \in \mathbf{O}_L \\ \theta_p \propto \frac{\|x_p x_{p_4}\| \cdot \|y_p y_{p_1}\|}{\|x_p x_{p_0}\| \cdot \|y_p y_{p_1}\|} & \text{如果 } \mathbf{P} \in \mathbf{O}_R \end{cases} \quad (2)$$

其中, $\mathbf{R}_{\text{Clo}}(\theta_p)$ 表示点 P 以唇区边界曲线为轴,向唇内旋转 θ_p 角的位移旋转运动,它与点 P 相对于各特征点的位置有关; $T_{\text{Clo}}(d_p)$ 为点 P 以唇区边界曲线为参考,位移为 d_p 的平移运动; x_p 为点 P 的横坐标; y_p 为点 P 的纵坐标; \mathbf{O}_L 为左半唇区域; \mathbf{O}_R 为右半唇区域。类似可得到下唇运动描述。

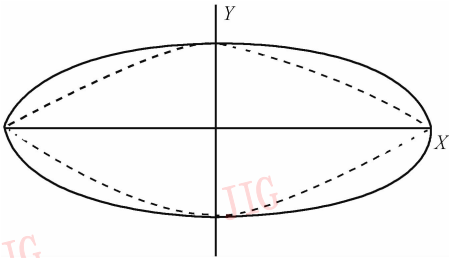


图4 突唇和翘唇运动旋转的曲线轴示意图

Fig.4 The curve axes of protruding and tilting lips

同样的模型也适用于内收嘴唇运动,仅仅在旋转和平移运动的幅度上与闭唇模型有差别。

与前述的闭唇运动类似,突唇运动也由旋转和平移构成,但在运动的参考点、方向和幅度上都有很大区别。在突唇运动中,口轮匝肌在唇区影响域内径向收缩,上下唇以上下唇区边界曲线(如图4椭圆中上下边界曲线所示)为轴,略向外转动,且向前突出。因此,上唇区域中某一点 P 的运动描述如下:

$$D_{\text{Pir}}(\mathbf{P}) = \mathbf{R}_{\text{Pir}}(\theta_p) + T_{\text{Pir}}(d_p) + C_{\text{Pir}}(a_p)$$

$$\begin{cases} \theta_p \propto \exp\left\{\frac{\|x_p x_{p_0}\|}{m}\right\} & \mathbf{P} \in \mathbf{O}_L \\ \theta_p \propto \exp\left\{\frac{-\|x_p x_{p_4}\|}{m}\right\} & \mathbf{P} \in \mathbf{O}_R \end{cases} \quad (3)$$

其中, $\mathbf{R}_{\text{Pir}}(\theta_p)$ 为点 P 以唇区边界曲线为轴,向唇外

侧旋转 θ_p 角的位移旋转运动; m 为旋转角 θ_p 的控制参数; $T_{\text{Pir}}(d_p)$ 为点 P 以唇区边界曲线为参考,垂直于唇区椭圆平面向外方向,位移为 d_p 的平移运动; $C_{\text{Pir}}(a_p)$ 为点 P 以唇区边界曲线为参考,径向位移为 a_p 的径向收缩运动,由 Waters 提出的肌肉模型得到; \mathbf{O}_L 为左半唇区域; \mathbf{O}_R 为右半唇区域。类似可得到下唇突唇运动描述。

同样的模型也适用于翘唇运动,但翘唇运动中的旋转子运动相对复杂,其旋转轴不再是唇区边界。如图3所示,以点 P_1 和点 P_2 的连线作为嘴唇的纵中轴线 Y ,点 P_0 和点 P_4 的连线作为嘴唇的横中轴线 X 。若唇区某点 P 越靠近 Y 轴,则该点的旋转轴越远离 X 轴,反之则越靠近 X 轴。此时,旋转运动的旋转轴如图4中的虚线所示。

4 实验结果

为验证本文提出的唇区肌肉模型的合成效果,利用一组人脸图像进行了合成口型实验。实验首先通过拍摄一组正在说话的人脸视频,同时选取其中静音时的某一帧照片为所要合成人脸的纹理信息,并将其与一个3维人脸模型相对应来合成一张3维人脸;然后,对唇区进行预处理,分割出上下唇;最后,添加牙齿与舌头的3维模型,以增加人脸的真实感。合成的无表情人脸如图5(a)所示,其具有2079个网格控制点和3919个多边形面片。

图5(b)~图5(e)分别为合成的闭唇、收唇、突唇与翘唇口型。其合成过程的描述如下:

首先,在人脸模型上标注如图3所示的5个特征点,并由此确定椭圆唇区;

然后,定义唇区各运动的正方向。如图6所示的3维坐标系中, X 、 Y 、 Z 轴的正方向分别为平移运动三分量的正方向,从 X 轴正方向看去,以顺时针方向为旋转运动的正方向。

接下来,调试并确定各运动参数。其中,既包括恒定参数,如旋转轴、平移参考坐标等,又包括可变参数,如旋转角、平移位移与 Waters 提出的模型中的一些肌肉参数。

在闭唇模型中,上唇的旋转运动 $\mathbf{R}_{\text{Clo}}^+(\theta_p)$ 中的旋转角为 0.15rad ,平移运动 $T_{\text{Clo}}^+(d_p)$ 中的点 p 的位移 $d_p = -0.05$ 单位距离;下唇收唇运动 $\mathbf{R}_{\text{Add}}^-(\theta_p)$ 中的旋转角为 -0.25rad ,平移运动 $T_{\text{Add}}^-(d_p)$ 中的位移

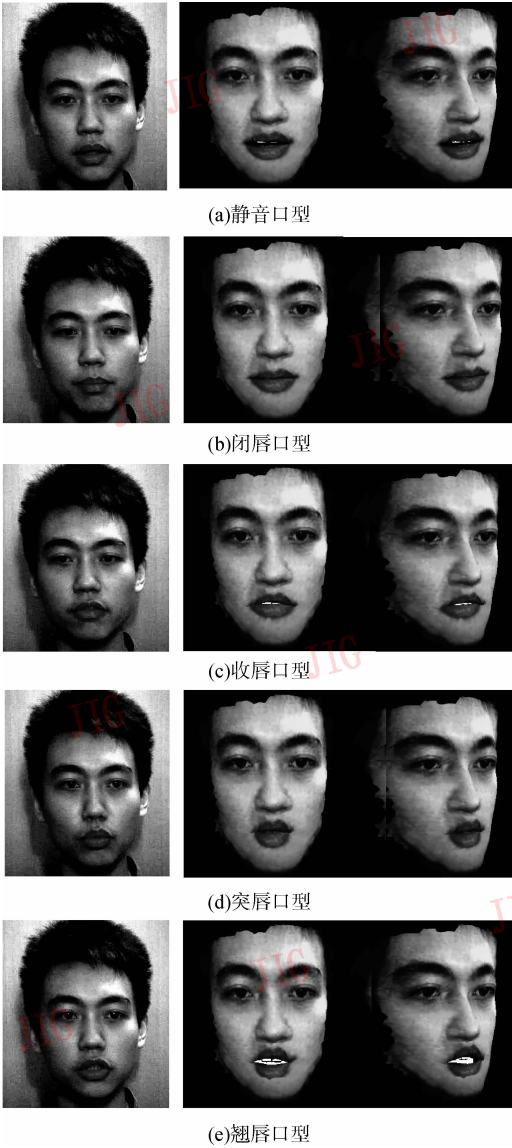


图 5 人脸及其口型的合成图像
Fig. 5 Face and synthesized lip shapes

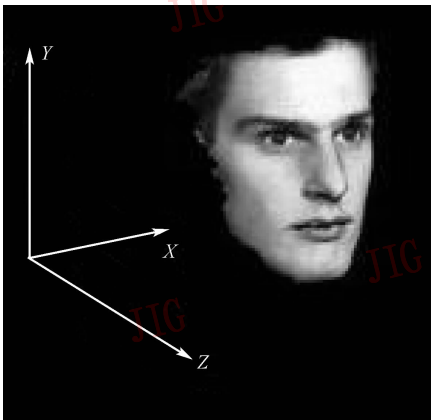
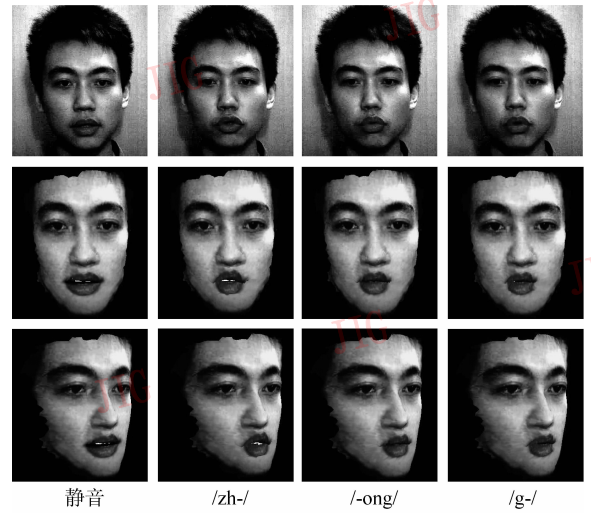


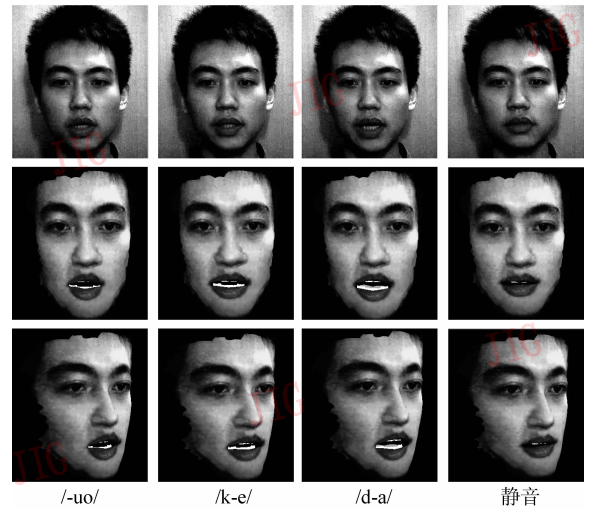
图 6 合成人脸的 3 维坐标系
Fig. 6 The 3D rectangular coordinates of synthesized face

d_p 为 0.06 单位距离 (合成系统中的距离单位); 突唇模型的口轮匝肌影响域中, 最大径向收缩位移 $a_p = 0.17$ 单位距离, 上唇旋转运动 $R_{\text{plr}}^+(\theta_p)$ 中的最大旋转角为 -0.2rad , 下唇最大旋转角为 0.2rad ; 翘唇运动中, 上唇旋转运动 $R_{\text{plr}}^+(\theta_p)$ 的最大旋转角为 -0.35rad , 下唇的最大旋转角为 0.35rad 。

图 7 为实拍的说话人在连续语流中的口型视频关键帧序列, 以及运用本文提出的模型所合成的对应关键帧口型。在拍摄过程中, 演员以约 2 字/s 的语速朗读词语“中国科大”。实拍图像为正面效果, 合成图像为正面与侧面效果。



(a) 第 1 组关键帧口型合成



(b) 第 2 组关键帧口型合成

图 7 对于连续语流“中国科大”合成的关键帧口型 (每一组上方是实拍视频关键帧, 中间是合成的关键帧口型, 下方是合成的侧面效果)

Fig. 7 The synthesized key lip shape sequence of reading “zhong guo ke da” (The upper row shows the real key lip shapes; The middle row shows the synthesized frontal ones; The lower row shows the lateral ones)

5 结 论

本文在 Waters 提出的肌肉模型的基础上,提出了一种新的唇区肌肉模型。利用该模型,并通过各肌肉参数的调节和不同肌肉群的组合运动合成了说话时的各种可能口型,并由这些合成的口型,通过帧间插值的方法生成了连续语流条件下的口型动画。实验结果表明,本研究合成的口型具有较高的自然度,从而弥补了 Waters 提出的肌肉模型在唇区建模上的不足。

本研究在许多方面都还可以进行拓展。例如在肌肉模型方面,虽然建立的模型可以合成普通说话时的各种口型,但所合成口型的幅度对于整体自然度的影响还是很大的。当口型幅度系数很大,使嘴唇动作非常夸张时,自然度会有一定的降低。虽然普通语音动画并不需要非常夸张的口型,但对于一个完善的嘴唇模型来说,应该要做到对任何可能的口型都能进行精确的模拟。由此可见,在整体模型的完善性上还需要做一些工作,例如首先对于夸张口型和复杂的非说话口型都能精确地进行描述;其次,由于每个人的面部结构、个人习惯、和种族文化的差异,在说话时所表达的口型也会有其个性化的一面,而该唇区肌肉模型却反映共性成分较多,个性成分相对较少,因此在今后的工作中,需要开展口型个性特征的抽取与控制方面拓展的研究。

另外,如何将该肌肉模型用于连续语流的语音动画合成,以便生成具有很高自然度的口型动画,也是一个可以拓展的方向。因此,今后拟尝试利用肌肉模型和与语音相关的面部解剖学知识,结合语音语言学的研究成果来建立语音与视觉表现上的动态关联规则,以进一步从语音语言学角度研究动态语

音动画的建模问题。

参考文献 (References)

- 1 Parke F I, Waters K. Computer Facial Animation [M]. Boston, MA, USA: Wellesley, 1996: 1 ~ 365.
- 2 Williams L. Performance driven facial animation [J]. Computer Graphics (ACM SIGGRAPH'90), 1990, 24(4): 235 ~ 242.
- 3 Terzopoulos D, Waters K. Physically-based facial modeling, analysis, and animation [J]. Journal of Visualization and Computer Animation, 1990, 1 (4): 73 ~ 80.
- 4 Waters K. A muscle model for animating three dimensional facial expression [J]. Computer Graphics (ACM SIGGRAPH'87), 1987, 22 (4): 17 ~ 24.
- 5 Blanz V, Vetter T. A morphable model for the synthesis of 3D faces [A]. In: Proceedings of ACM SIGGRAPH '99 Conference Proceedings [C], Los Angeles, CA, USA, 1999: 187 ~ 194.
- 6 Albrecht I, Haber J, Kahler K, et al. 'May I talk to you? :-)'—Facial animation from text [A]. In: Proceedings of Pacific Graphics [C], Beijing, China, 2002: 77 ~ 86.
- 7 Kshirsagar S, Magnenat-Thalmann N. Lip synchronization using linear predictive analysis [A]. In: Proceedings of IEEE International Conference on Multimedia and Expo [C], New York, USA, 2000: 1077 ~ 1080.
- 8 Song Ming-li, Chen Chun, Bu Jia-jun, et al. 3D realistic talking face co-driven by text and speech [A]. In: Proceedings of IEEE International Conference on Systems, Man and Cybernetics [C], Washington, DC, USA, 2003: 2175 ~ 2186.
- 9 Pi Xin, et al. Oral Cavity Anatomy [M]. Beijing: People's Medical Publishing House, 1987: 111 ~ 115. [皮昕等. 口腔解剖生理学 [M], 北京: 人民卫生出版社, 1987: 111 ~ 115.]
- 10 Epstein M, Hacopian N, Ladefoge P. Dissection of the Speech Production Mechanism [M]. Los Angeles, CA, USA: The UCLA Phonetics Laboratory, 2002: 12 ~ 15.
- 11 Moubaraki L, Ohya J. Realistic 3D mouth animation using a minimal number of parameters [A], In: Proceedings of IEEE International Workshop on Robot and Human Communication [C], Tsukuba, Japan, 1996: 201 ~ 206.