

SISTEMA DE ANIMAÇÃO FACIAL TRIDIMENSIONAL E SÍNTESE DE VOZ



Aluna: Andréa Britto Mattos

Orientador: Roberto Marcondes Cesar Jr.

INTRODUÇÃO

Este trabalho tem como objetivo complementar o projeto *Avatar*, que teve início em 2006 pelo aluno Marcos Moreti[1], continuando em 2007, com a aluna Flávia Ost[2]. Nele, o usuário poderia digitar perguntas para uma face tridimensional que pronunciava as respostas. Além disso, era usada uma câmera para que a face apenas respondesse quando detectava movimento, o que foi feito utilizando a biblioteca OpenCV, de processamento de vídeo.

O Avatar utilizava o programa Haptik para cuidar da animação, síntese de voz e sincronização labial do modelo. No entanto, o Haptik impunha muitas restrições para o sistema: o Avatar deveria usar o Windows como plataforma, o Visual Studio como IDE e a biblioteca MFC (*Microsoft Foundation Classes*), não disponível na versão gratuita do Visual Studio.

Assim, a proposta deste projeto é a substituição do Haptik por um sistema próprio de animação e síntese de voz, utilizando um ambiente gratuito.

INTERPOLAÇÃO DE QUADROS-CHAVES

O sistema de animação utilizou a técnica de *Morph Target*, na qual um conjunto fixo de quadros-chaves são interpolados pelo computador.

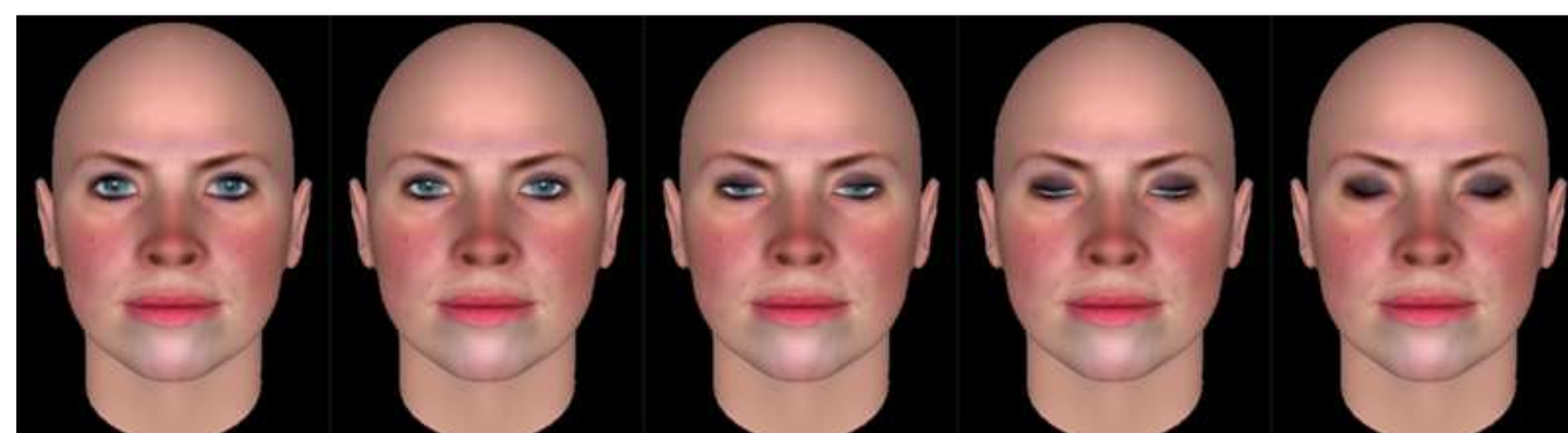


FIGURA 1: Dadas duas poses-chaves - face de olhos abertos e de olhos fechados - o processo de interpolação gera estados intermediários, evitando transições bruscas

As poses exibidas no programa são divididas em dois grupos: *visemas* e *expressões faciais*. Um visema é o correspondente visual de um fonema, podendo corresponder a mais de um fonema. Para as expressões faciais, foi utilizada a proposta de Paul Ekman[4], na qual seis emoções puras (alegria, tristeza, raiva, medo, nojo e surpresa) podem gerar todas as outras.

FERRAMENTAS

Para a modelar as faces e exportá-las para o sistema, foi utilizado o programa FaceGen, capaz de gerar expressões faciais, visemas e alguns movimentos não-verbais. Para a síntese de voz, foi utilizado o programa eSpeak, capaz de extrair os fonemas de um texto em português, e pronunciá-lo.

Para a animação, foi usada a engine gráfica Ogre 3D, que permite interpolar quadros-chaves e combinar *poses* com diferentes *influências*. Uma pose constitui um conjunto de valores que definem a deformação de cada vértice da malha tridimensional. A influência é a quantidade de deformação da pose. Para obter uma animação suave, foi preciso incrementar, em cada frame, a influência das poses exibidas.

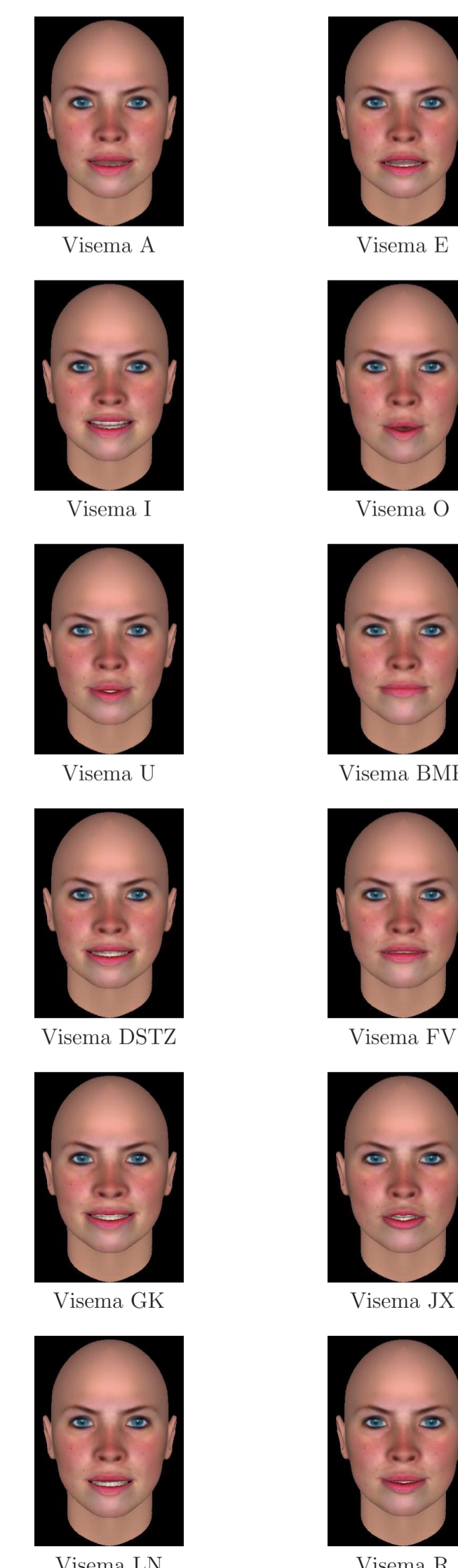
IMPLEMENTAÇÃO

A importação dos modelos do FaceGen para o Ogre teve que respeitar as restrições de formato e animação da engine. A conversão desses modelos foi feita em dois passos: primeiro, os modelos foram convertidos para uma especificação de XML. Em seguida, foi implementado um script para modificar diretamente estes arquivos XML.

Para a sincronização labial, foi implementada uma callback que é chamada automaticamente pelo eSpeak, sempre que um buffer de áudio é produzido. Esta callback pode retornar eventos de fonema na medida em que ocorrem na fala. Estes fonemas são mapeados para o visema correspondente, que é, por fim, exibido. O mapeamento adotado é mostrado a seguir.

Mapeamento de Fonemas para Visemas			
Fonema no IPA ¹	Fonema em ASCII	Exemplo do som	Visema Correspondente
/a/	a	saco	Visema A
/e/	e	seco	Visema E
/ɛ/	E	bela	
/i/	i	sico	Visema I
/o/	o	soco	Visema O
/ɔ/	O	sozinho	
/u/	u	suco	Visema U
/b/	b	bata	
/m/	m	mata	Visema BMP
/p/	p	pata	
/d/	d	data	
/s/	s	saca	Visema DSTZ
/t/	t	tata	
/z/	z	zaca	
/f/	f	faca	Visema FV
/v/	v	vaca	
/g/	g	gata	Visema GK
/k/	k	kata	
/ʒ/	Z	jaca	Visema JX
/j/	S	chaga	
/l/	l	galo	
/ʎ/	L	galho	Visema LN
/n/	n	nata	
/ɲ/	N	ganho	
/ɣ/	G	carro	Visema R
/r/	R	caro	

TABELA 1: Mapeamento de fonemas do português - no IPA e em ASCII - para grupos de visemas, exibidos nas imagens ao lado.



¹ International Phonetic Alphabet

VISÃO COMPUTACIONAL

Nesta versão, o algoritmo anterior de detecção de movimento foi substituído pela implementação do OpenCV de reconhecimento facial, para que o avatar exigisse a presença de um interlocutor humano. Foi adicionada também outra funcionalidade: mover os olhos da face conforme a posição do usuário.

RESULTADOS

Não há Inteligência Artificial no módulo de animação, de forma que o avatar apenas repete as frases digitadas. Abaixo estão algumas *screenshots* da tela do programa. Vídeos foram também adicionados na página do projeto[3].

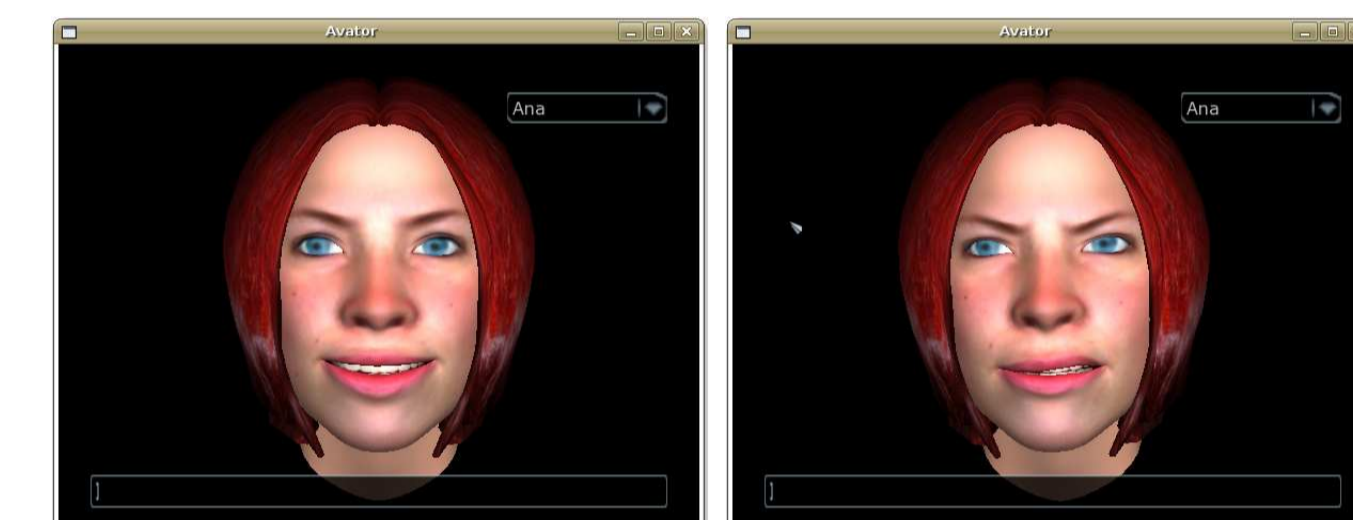


FIGURA 2: Avatar exibindo expressões faciais de alegria (à esquerda) e nojo (à direita)



FIGURA 3: Diferentes modelos que podem ser carregados no programa

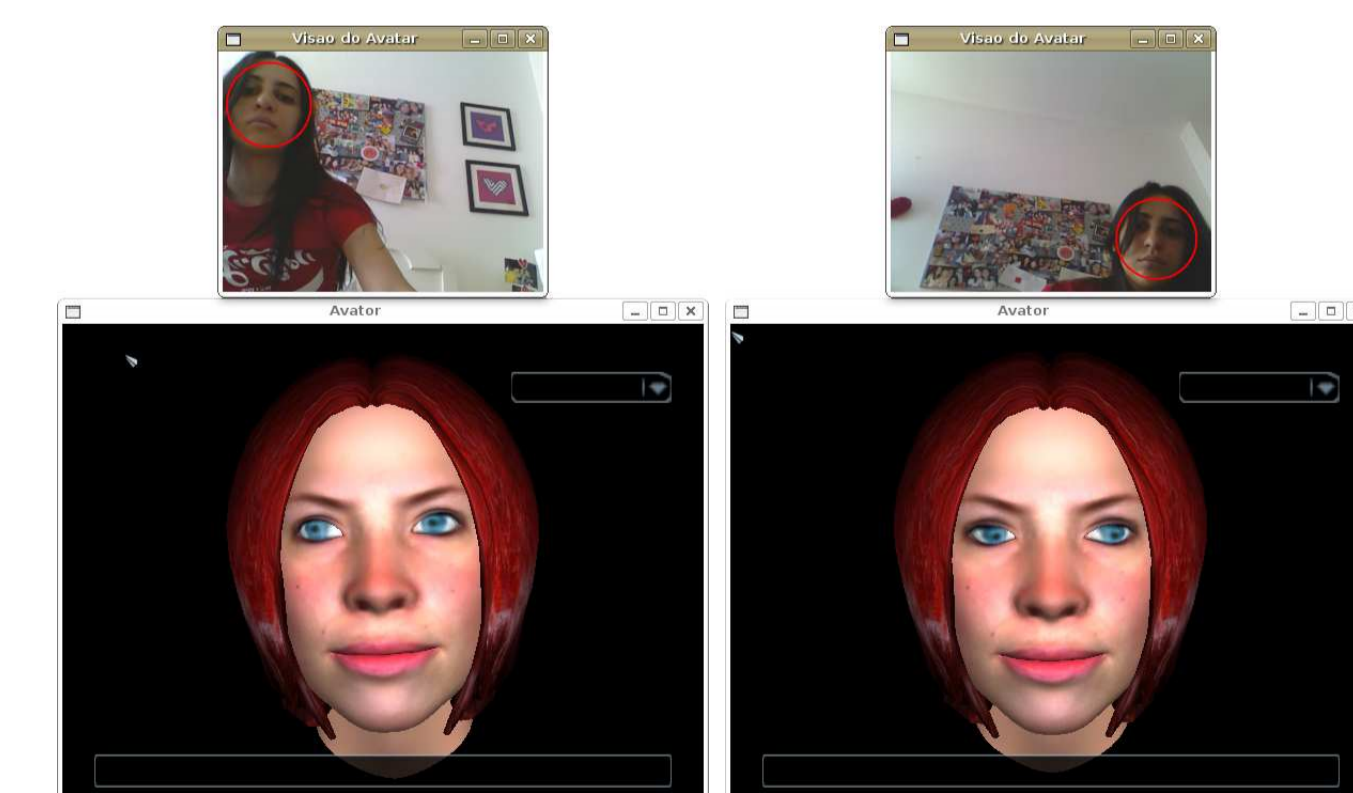


FIGURA 4: Movimento dos olhos do avatar seguindo a face encontrada

Referências

- [1] <http://www.linux.ime.usp.br/~cef/mac499-06/monografias/mpmoreti/>
- [2] <http://www.vision.ime.usp.br/~fost/>
- [3] <http://www.linux.ime.usp.br/~dedea/mac499/>
- [4] EKMAN, P. All Emotions are Basic. Oxford University Press, Nova York, 1994.