

# OpenAI 推出首个文生视频大模型 Sora，引领 AI 文生视频行业跨越式发展

——计算机行业跟踪报告

强于大市 (维持)

2024 年 02 月 18 日

## 行业核心观点:

文生视频大模型 Sora 重磅发布，可生成长达 1 分钟的视频。2 月 16 日，OpenAI 推出其首个文生视频大模型 Sora。根据官网介绍，Sora 可以生成长达 1 分钟时长的视频，同时还能保证视频质量，并遵循用户的提示 (prompt)。

## 投资要点:

Sora 是一个扩散 transformer，具有强大的语言理解能力，通过在潜在空间训练 patches 生成视频。对标 tokens，OpenAI 将视觉数据转换为 patches，有效用于 Sora 大模型训练。Sora 是一种扩散模型，通过给出输入的静态噪声以及相关的文本提示 (prompt) 等调节信息，训练生成原始的“干净”patches。在推理时，OpenAI 还可以通过在适当大小的网格中排列随机初始化的 patches 来控制生成视频的大小。与 GPT 模型类似，Sora 使用 transformer 架构，释放出卓越的扩展性能。立足 DALL·E 3 和 GPT 模型，Sora 具有强大的语言理解能力，能够生成更加准确遵循用户提示的高质量视频。此外，在固定种子和输入的情况下，可以看到训练计算的增加能显著提升样本视频的质量。

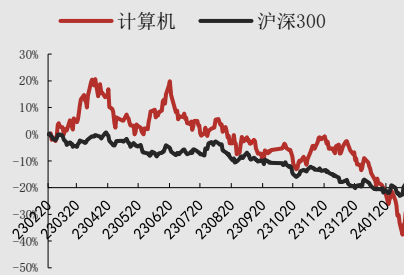
**多维度跨越式突破，视频质量飞跃性提升。**Sora 的采样更具有灵活性，同时改进了框架和构图。Sora 可以采样宽屏 1920x1080p 的视频、垂直 1080x1920 的视频以及介于两者之间的所有视频。这让 Sora 可直接以不同的原始长宽比创建内容。OpenAI 还通过经验发现，在视频的原始长宽比上进行训练可以改善构图和框架。Sora 还支持图生视频、视频生视频，能执行广泛的图像和视频编辑任务，创建完美的循环视频、动画静态图像、向前或向后扩展视频等。在连接视频上，Sora 能将两个输入视频无缝衔接在一起。虽然目前 Sora 仍然有一些缺陷和局限性，但已经开始理解物理意义，并出现许多有趣的涌现能力，如三维一致性。

**重塑 AI 文生视频行业格局，或冲击 AI 文生图赛道。**Sora 在生成视频长度上大幅领先，多角度镜头能力也显著领先行业竞品。同样的 prompt，Sora 生成的视频长度、质量都显著领先。Sora 可以生成可变大小的图像，最高可达 2048 × 2048 分辨率，图片画质有了大幅提升。我们认为随着文生视频画质能力的提升，图片作为单帧的视频，文生视频领域的产品或将冲击文生图行业。

**投资建议:** 1) AI 文生视频行业发展带动 AI 行业应用落地的机遇; 2) AI 行业发展对算力、光模块等基础设施的持续需求; 3) AIGC 在媒体、游戏等行业的加速落地带来的投资机遇。

**风险提示:** AI 产业发展不及预期; AI 带来的版权、隐私及技术风险; 国内 AI 应用落地不及预期; 中美科技摩擦风险。

## 行业相对沪深 300 指数表现



数据来源: 聚源, 万联证券研究所

## 相关研究

Q4 基金重仓略微超配, 前十大重仓股组成不变

人工智能行业应用多点开花

工信部就《国家人工智能产业综合标准化体系建设指南》公开征求意见

分析师:

夏清莹

执业证书编号:

S0270520050001

电话:

075583223620

邮箱:

xiaqy1@wlzq.com.cn

## 正文目录

<b>1 OpenAI 发布 Sora, AI 文生视频大模型跨越性突破</b> .....	<b>3</b>
1.1 OpenAI 首个文生视频大模型 SORA 重磅推出.....	3
1.2 多维度跨越式突破, 视频质量飞跃性提升.....	5
1.3 重塑 AI 文生视频行业格局, 或冲击 AI 文生图赛道.....	7
<b>2 投资建议</b> .....	<b>9</b>
<b>3 风险提示</b> .....	<b>9</b>
图表 1: Sora 一分钟展示视频的 prompt 及部分截图.....	3
图表 2: Sora 将视觉数据转换为 patches 的示意图.....	3
图表 3: Sora 通过扩散还原视频的示意图.....	4
图表 4: 不同训练计算生成的样本视频对比.....	4
图表 5: 使用正方形裁剪(左)与使用原始大小(右)的训练视频效果对比.....	5
图表 6: 向后扩展视频示意.....	5
图表 7: 从左上图逐渐转化至右下图的场景示意.....	6
图表 8: Sora 三维一致性示意图.....	6
图表 9: 其他文生视频产品的部分参数统计.....	7
图表 10: 相同 prompt 的生成视频成果对比.....	8
图表 11: Sora 的图像生成样本.....	8

# 1 OpenAI 发布 Sora, AI 文生视频大模型跨越性突破

## 1.1 OpenAI 首个文生视频大模型 SORA 重磅推出

文生视频大模型Sora重磅发布,可生成长达1分钟的视频。2月16日,OpenAI推出其首个文生视频大模型Sora。根据官网介绍,Sora可以生成长达1分钟时长的视频,同时还能保证视频质量,并遵循用户的提示(prompt)。

图表1: Sora 一分钟展示视频的 prompt 及部分截图

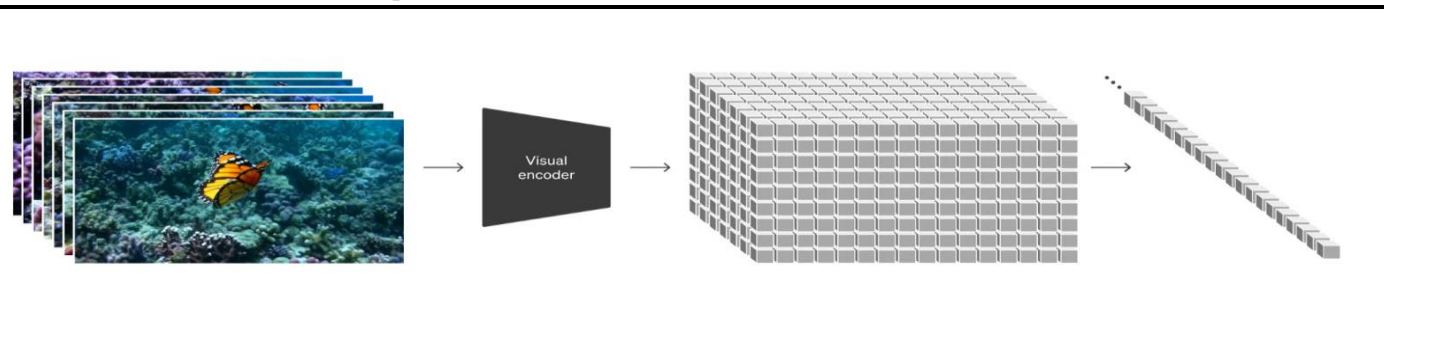
Prompt (提示)	视频部分截图
<p>A stylish woman walks down a Tokyo street filled with warm glowing neon and animated city signage. She wears a black leather jacket, a long red dress, and black boots, and carries a black purse. She wears sunglasses and red lipstick. She walks confidently and casually. The street is damp and reflective, creating a mirror effect of the colorful lights. Many pedestrians walk about.</p> <p>翻译: 一位时尚的女人走在东京的街道上,街道上到处都是温暖的发光霓虹灯和动画城市标志。她身穿黑色皮夹克,红色长裙,黑色靴子,背着一个黑色钱包。她戴着墨镜,涂着红色口红。她自信而随意地走路。街道潮湿而反光,营造出五颜六色的灯光的镜面效果。许多行人四处走动。</p>	

资料来源: OpenAI, 万联证券研究所

注: 翻译内容来自Microsoft Edge网页自带翻译。

将视觉数据转换为patches,有效用于Sora大模型训练。LLM范式的成功部分受益于使用tokens,tokens能够将文本的多种模态(代码、数学、各种自然语言)统一起来。OpenAI基于LLMs使用文本tokens的灵感,将所有视觉数据转化为patches,在Sora中实现类似的效果。根据OpenAI的介绍,patches此前就已经被证明是视觉数据模型的有效表示,同时OpenAI还发现,patches在训练生成不同类型视频和图像模型中是一种高度可扩展且有效的表示。

图表2: Sora 将视觉数据转换为 patches 的示意图

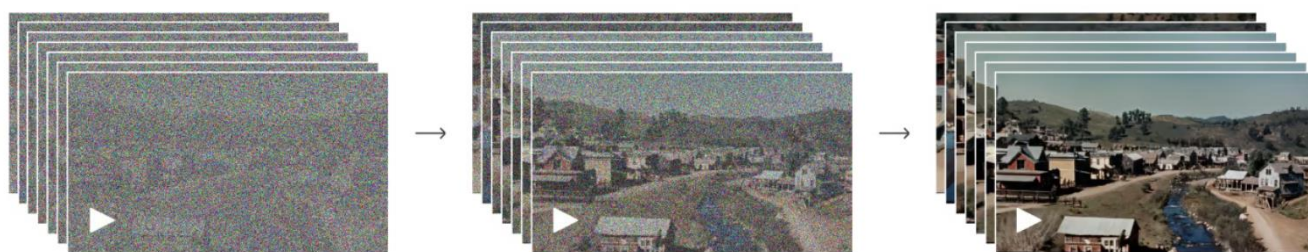


资料来源: OpenAI, 万联证券研究所



Sora是一个扩散transformer (diffusion transformer)，通过在潜在空间训练patches生成视频。具体来看视频生成的过程，1)首先将视频压缩到低维的潜在空间：OpenAI训练了一个降低视觉数据维度的网络，通过这个网络原始视频会在时间和空间上都被压缩，并输出为潜在表示；2)用时空潜在patches训练Sora：Sora在这个压缩后的潜在空间中接受训练，基于从原始视频中提取的时空潜在patches，OpenAI能够使得Sora对不同分辨率、持续时间和长宽比的视频和图像进行训练（图像相当于单帧视频）；3)解码生成新视频：OpenAI训练了对应的解码器模型，将Sora在潜在空间训练生成的视频（潜在表示）映射回像素空间；在推理时，OpenAI还可以通过在适当大小的网格中排列随机初始化的patches来控制生成视频的大小。**Sora是一种扩散模型**，通过给出输入的静态噪声以及相关的文本提示（prompt）等调节信息，训练生成原始的“干净”patches。与GPT模型类似，Sora使用transformer架构，释放出卓越的扩展性能。

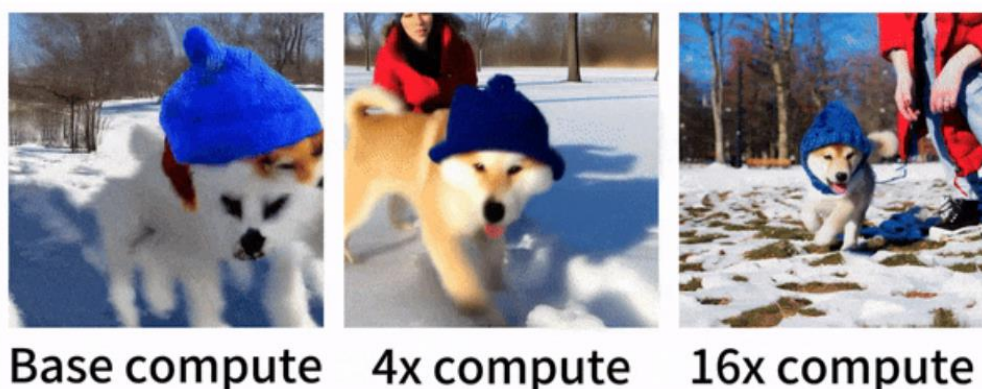
图表3: Sora 通过扩散还原视频的示意图



资料来源: OpenAI, 万联证券研究所

**训练计算的增加可以显著提升视频质量。**在固定种子和输入的情况下，可以看到训练计算的增加能显著提升样本视频的质量。

图表4: 不同训练计算生成的样本视频对比



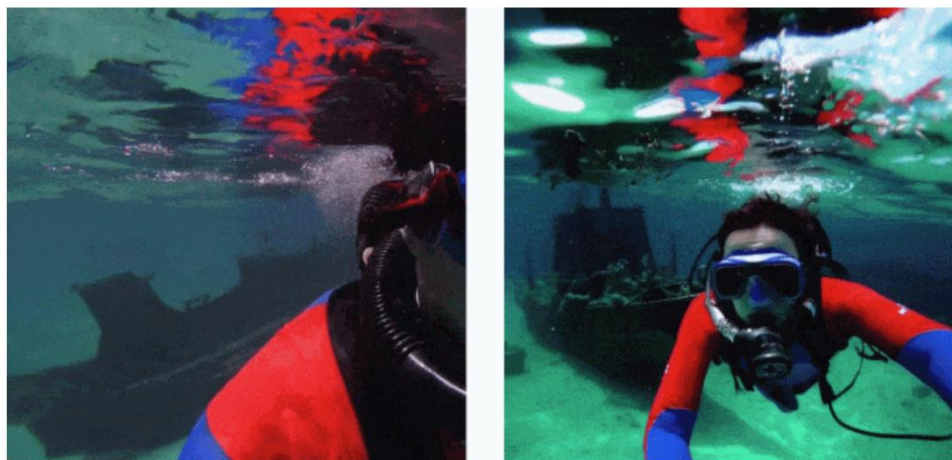
资料来源: OpenAI, 万联证券研究所

**立足DALL·E 3和GPT模型，Sora具有强大的语言理解能力，能够生成更准确、更高质量的视频。**OpenAI将DALL·E 3的re-captioning技术应用在Sora中，先训练一个高度描述性的字幕生成器模型，然后使用它为训练集中所有视频生成文本字幕，通过对高度描述性视频字幕进行训练，能够提供文本的保真度以及视频的整体质量。此外，OpenAI还利用GPT将简短的用户prompt转化为较长的详细字幕，然后发送到视频模型中，使得Sora能够生成更加准确遵循用户提示的高质量视频。

## 1.2 多维度跨越式突破，视频质量飞跃性提升

Sora的采样更具有灵活性，同时改进了框架和构图。过去的图像和视频生成方法通常需要调整大小、进行裁剪或者是将视频剪切到标准尺寸，例如4秒的视频分辨率为256x256。而OpenAI的研究发现在原始大小的数据上进行训练，采样更具灵活性、同时可以提高视频质量。Sora可以采样宽屏1920x1080p的视频、垂直1080x1920的视频以及介于两者之间的所有视频。这让Sora可直接以不同的原始长宽比创建内容。Sora还支持在生成全分辨率的内容之前，以较小的尺寸快速创建内容原型——所有内容都使用相同的模型。OpenAI还通过经验发现，在视频的原始长宽比上进行训练可以改善构图和框架。研究团队将Sora与其他模型的一个版本进行比较，该版本将所有训练视频裁剪为方形。在正方形裁剪上训练的模型有时会生成仅部分显示主题的视频（左图），相比之下，来自Sora的视频有改进的帧内容（右图）。

图表5: 使用正方形裁剪（左）与使用原始大小（右）的训练视频效果对比



资料来源: OpenAI, 万联证券研究所

Sora还支持图生视频、视频生视频，能够向前/向后扩展视频。Sora能基于DALL·E 3的图像生成视频，还能执行广泛的图像和视频编辑任务，创建完美的循环视频、动画静态图像、向前或向后扩展视频等。以下是Sora从一段生成的视频向后拓展出的三个新视频，可以看到新视频的开头各不相同，但拥有相同的结尾。

图表6: 向后扩展视频示意



资料来源: OpenAI, 万联证券研究所



在连接视频上，Sora能将两个输入视频无缝衔接在一起。如下图所示，左上图是村庄，右下图是海洋，看似毫不相关的场景，Sora通过逐渐放大河流，合理、无缝地从左上→左下→右上→右下实现两个视频的转化连接。

图表7: 从左上图逐渐转化至右下图的场景示意



资料来源: OpenAI, 万联证券研究所

虽然目前Sora仍然有一些缺陷和局限性，但已经开始理解物理世界，并出现许多有趣的涌现能力。Sora目前还存在许多局限性，例如不能准确模拟许多基本交互的物理现象，如玻璃碎裂；如吃食物，并不总能产生正确的物体状态变化。但我们认为Sora已经接触到了世界模型的范畴。Sora能够生成具有多个角色、特定类型的运动以及主题和背景准确细节的复杂场景。该模型不仅了解用户在提示中要求的内容，还了解这些东西在物理世界中的存在方式。OpenAI发现，视频模型在经过大规模训练后，会表现出许多有趣的新能力。这些能力使Sora能够模拟物理世界中的人、动物和环境的某些方面。这些特性的出现没有任何明确的三维、物体等归纳偏差，它们纯粹是规模现象。例如三维一致性：Sora可以生成动态摄像机运动的视频，随着摄像机的移动和旋转，人物和场景元素在三维空间中的移动是一致的。

图表8: Sora 三维一致性示意图



资料来源: OpenAI, 万联证券研究所

### 1.3 重塑 AI 文生视频行业格局，或冲击 AI 文生图赛道

Sora在生成视频长度上大幅领先，多角度镜头能力也显著领先行业竞品。与此前的文生视频产品相比，Sora在视频长度上实现了质的飞跃，此前的产品还停留在4s左右拓展的阶段，拓展的连贯性一直是众多产品的研发重点。而此次Sora则是直接能够生成最长60s的视频，这对AI文生视频行业来说是质的突破。同时，此前的文生视频产品只能通过单角度镜头生成视频，而Sora可以在单个生成的视频中创建多个角度的镜头，以准确保留角色和视觉风格，更准确地解释提示并生成表达生动情感的引人注目的角色。

图表9: 其他文生视频产品的部分参数统计

Company	Generation Type	Max Length	Extend?	Camera Controls? (zoom, pan)	Motion Control? (amount)	Other Features	Format
Runway	Text-to-video, image-to-video, video-to-video	4 sec	Yes	Yes	Yes	Motion brush, upscale	Website
Pika	Text-to-video, image-to-video	3 sec	Yes	Yes	Yes	Modify region, expand canvas, upscale	Website
Genmo	Text-to-video, image-to-video	6 sec	No	Yes	Yes	FX presets	Website
Kaiber	Text-to-video, image-to-video, video-to-video	16 sec	No	No	No	Sync to music	Website
Stability	Image-to-video	4 sec	No	No	Yes		Local model, SDK
Zeroscope	Text-to-video	3 sec	No	No	No		Local model
ModelScope	Text-to-video	3 sec	No	No	No		Local model
AnimateDiff	Text-to-video, image-to-video, video-to-video	3 sec	No	No	No		Local model
Morph	Text-to-video	3 sec	No	No	No		Discord bot
Hotshot	Text-to-video	2 sec	No	No	No		Website
Moonvalley	Text-to-video, image-to-video	3 sec	No	No	No		Discord bot
Deforum	Text-to-video	14 sec	No	Yes	No	FX presets	Discord bot
Leonardo	Image-to-video	4 sec	No	No	Yes		Website
Assistive	Text-to-video, image-to-video	4 sec	No	No	Yes		Website
Neural Frames	Text-to-video, image-to-video, video-to-video	Unlimited	No	No	No	Sync to music	Website
Magic Hour	Text-to-video, image-to-video, video-to-video	Unlimited	No	No	No	Face swap, sync to music	Website
Vispunk	Text-to-video	3 sec	No	Yes	No		Website
Decohere	Text-to-video, image-to-video	4 sec	No	No	Yes		Website
Domo AI	Image-to-video, video-to-video	3 sec	No	No	Yes		Discord bot
FullJourney	Text-to-video, image-to-video	8 sec	No	Yes	No	Lipsyncing, face swap	Discord bot

Source: @venturetwins  
Charts provided herein are for informational purposes only and should not be relied upon when making any investment decision. Past performance is not indicative of future results. None of the above should be taken as investment advice; please see a16z.com/disclosures for more information.

a16z Consumer

资料来源: a16z, 万联证券研究所

同样的prompt, Sora生成的视频长度、质量都显著领先。根据用户Gabor Cselle的测试, 将Sora和Pika、RunwayML和Stable Video进行了对比。在采用了与OpenAI示例中相同的prompt的基础上, 结果显示其他主流工具生成的视频都大约只有5秒钟, 而Sora可以生成一段长达17秒的视频, 并在视频场景中保持动作和画面一致性。



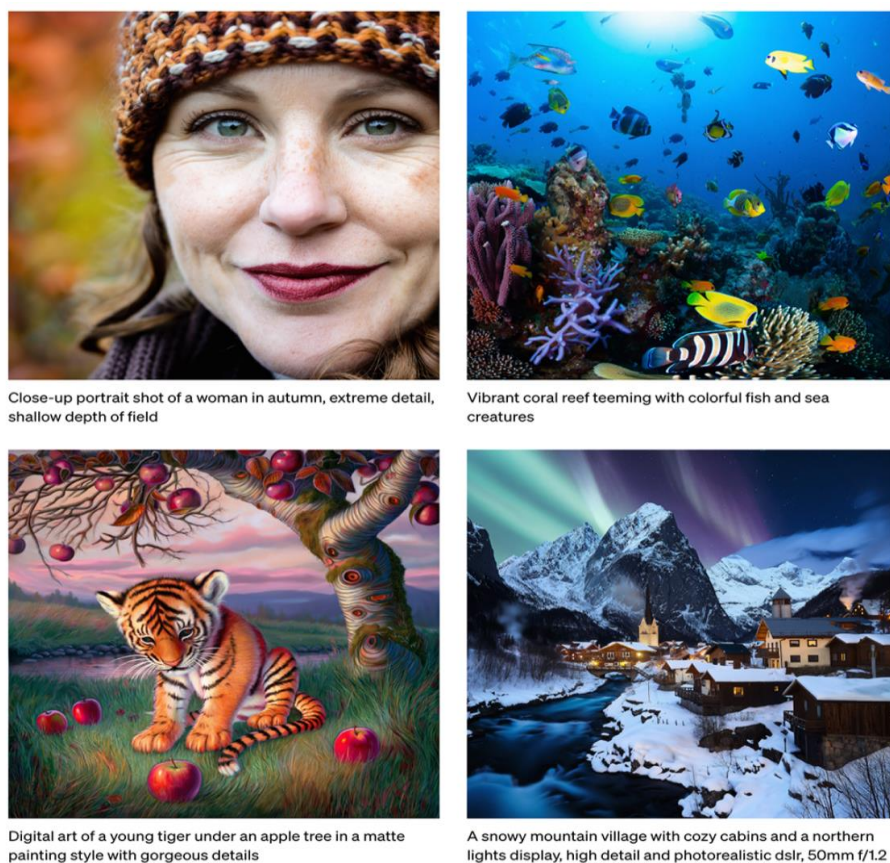
图表10: 相同 prompt 的生成视频成果对比



资料来源: 创业邦, 万联证券研究所

画质大幅提升, 或冲击AI文生图片行业。OpenAI通过在一个时间范围为一帧的空间网格中排列高斯噪声块来实现生成图片。该模型可以生成可变大小的图像, 最高可达 $2048 \times 2048$ 分辨率。可以看到, Sora生成的图片画质有了大幅提升, 非常清晰, 此前的文生视频产品由于清晰度不足的原因, 与文生图赛道暂未有强烈地竞争, 但我们认为随着文生视频画质能力的提升, 图片作为单帧的视频, 文生视频领域的产品或将冲击文生图行业。

图表11: Sora 的图像生成样本



资料来源: OpenAI, 万联证券研究所



## 2 投资建议

OpenAI推出的AI文生视频产品Sora，让文生视频行业实现了跨越式的发展，生成视频时长的突破、多角度镜头、强大的语言理解能力等都让Sora大幅领先行业竞品，也让AI文生视频、文生图、图生视频、视频生视频的落地性大幅提升，有望在多个行业领域应用落地。此外OpenAI的研究证明，训练计算的提升可以提升生成视频的样本质量，因此Sora的推出也展现出文生视频行业对大量算力的强烈需求。**建议关注：1) AI文生视频行业发展带动AI行业应用落地的机遇；2) AI行业发展对算力、光模块等基础设施的持续需求；3) AIGC在媒体、游戏等行业的加速落地带来的投资机遇。**

## 3 风险提示

AI产业发展不及预期；AI带来的版权、隐私及技术风险；国内AI应用落地不及预期；中美科技摩擦风险。

## 行业投资评级

强于大市：未来6个月内行业指数相对大盘涨幅10%以上；

同步大市：未来6个月内行业指数相对大盘涨幅10%至-10%之间；

弱于大市：未来6个月内行业指数相对大盘跌幅10%以上。

## 公司投资评级

买入：未来6个月内公司相对大盘涨幅15%以上；

增持：未来6个月内公司相对大盘涨幅5%至15%；

观望：未来6个月内公司相对大盘涨幅-5%至5%；

卖出：未来6个月内公司相对大盘跌幅5%以上。

基准指数：沪深300指数

## 风险提示

我们在此提醒您，不同证券研究机构采用不同的评级术语及评级标准。我们采用的是相对评级体系，表示投资的相对比重建议；投资者买入或者卖出证券的决定取决于个人的实际情况，比如当前的持仓结构以及其他需要考虑的因素。投资者应阅读整篇报告，以获取比较完整的观点与信息，不应仅仅依靠投资评级来推断结论。

## 证券分析师承诺

本人具有中国证券业协会授予的证券投资咨询执业资格并登记为证券分析师，以勤勉的执业态度，独立、客观地出具本报告。本报告清晰准确地反映了本人的研究观点。本人不曾因，不因，也将不会因本报告中的具体推荐意见或观点而直接或间接收到任何形式的补偿。

## 免责声明

万联证券股份有限公司（以下简称“本公司”）是一家覆盖证券经纪、投资银行、投资管理和证券咨询等多项业务的全国性综合类证券公司。本公司具有中国证监会许可的证券投资咨询业务资格。

本报告仅供本公司的客户使用。本公司不会因接收人收到本报告而视其为客户。在任何情况下，本报告中的信息或所表述的意见并不构成对任何人的投资建议。本报告中的信息或所表述的意见并未考虑到个别投资者的具体投资目的、财务状况以及特定需求。客户应自主作出投资决策并自行承担投资风险。本公司不对任何人因使用本报告中的内容所导致的损失负任何责任。在法律许可情况下，本公司或其关联机构可能会持有报告中提到的公司所发行的证券头寸并进行交易，还可能为这些公司提供或争取提供投资银行、财务顾问或类似的金融服务。

市场有风险，投资需谨慎。本报告是基于本公司认为可靠且已公开的信息撰写，本公司力求但不保证这些信息的准确性及完整性，也不保证文中的观点或陈述不会发生任何变更。在不同时期，本公司可发出与本报告所载资料、意见及推测不一致的报告。分析师任何形式的分享证券投资收益或者分担证券投资损失的书面或口头承诺均为无效。

本报告的版权仅为本公司所有，未经书面许可任何机构和个人不得以任何形式翻版、复制、刊登、发表和引用。未经我方许可而引用、刊发或转载的引起法律后果和造成我公司经济损失的概由对方承担，我公司保留追究的权利。

## 万联证券股份有限公司 研究所

上海浦东新区世纪大道 1528 号陆家嘴基金大厦

北京西城区平安里西大街 28 号中海国际中心

深圳福田区深南大道 2007 号金地中心

广州天河区珠江东路 11 号高德置地广场